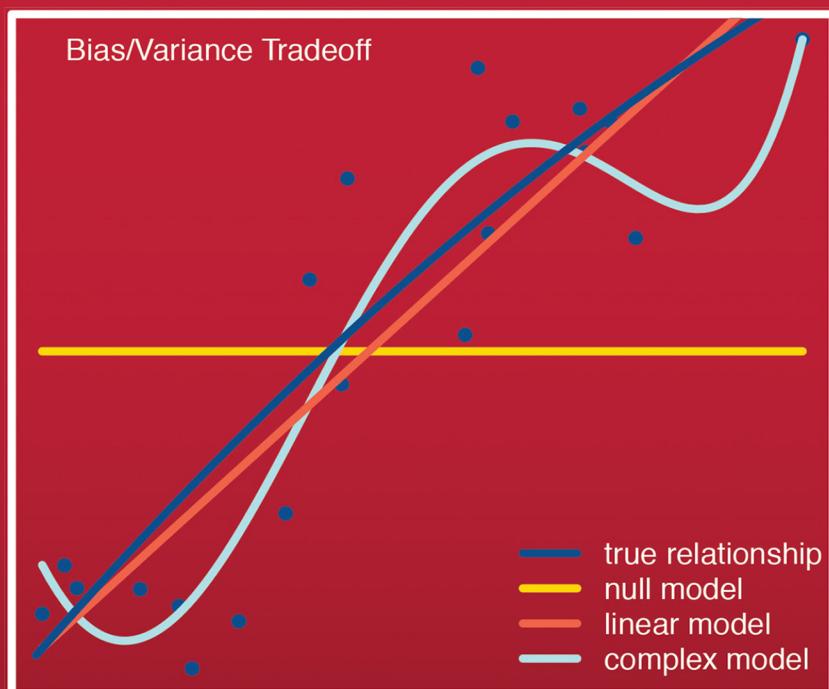


Texts in Statistical Science

Foundations of Statistics for Data Scientists

With R and Python



Alan Agresti

Maria Kateri



CRC Press
Taylor & Francis Group

A CHAPMAN & HALL BOOK

Foundations of Statistics for Data Scientists

CHAPMAN & HALL/CRC

Texts in Statistical Science Series

Joseph K. Blitzstein, *Harvard University, USA*

Julian J. Faraway, *University of Bath, UK*

Martin Tanner, *Northwestern University, USA*

Jim Zidek, *University of British Columbia, Canada*

Recently Published Titles

Randomization, Bootstrap and Monte Carlo Methods in Biology

Fourth Edition

Bryan F. J. Manly, Jorge A. Navarro Alberto

Principles of Uncertainty, Second Edition

Joseph B. Kadane

Beyond Multiple Linear Regression

Applied Generalized Linear Models and Multilevel Models in R

Paul Roback, Julie Legler

Bayesian Thinking in Biostatistics

Gary L. Rosner, Purushottam W. Laud, and Wesley O. Johnson

Linear Models with Python

Julian J. Faraway

Modern Data Science with R, Second Edition

Benjamin S. Baumer, Daniel T. Kaplan, and Nicholas J. Horton

Probability and Statistical Inference

From Basic Principles to Advanced Models

Miltiadis Mavrakakis and Jeremy Penzer

Bayesian Networks

With Examples in R, Second Edition

Marco Scutari and Jean-Baptiste Denis

Time Series

Modeling, Computation, and Inference, Second Edition

Raquel Prado, Marco A. R. Ferreira and Mike West

A First Course in Linear Model Theory

Second Edition

Nalini Ravishanker, Zhiyi Chi, Dipak K. Dey

Foundations of Statistics for Data Scientists

With R and Python

Alan Agresti and Maria Kateri

For more information about this series, please visit: <https://www.crcpress.com/Chapman--Hall/CRC-Texts-in-Statistical-Science/book-series/CHTEXSTASCI>

Foundations of Statistics for Data Scientists

With R and Python

Alan Agresti and Maria Kateri



CRC Press

Taylor & Francis Group

Boca Raton London New York

CRC Press is an imprint of the
Taylor & Francis Group, an **informa** business
A CHAPMAN & HALL BOOK

First edition published 2022
by CRC Press
6000 Broken Sound Parkway NW, Suite 300, Boca Raton, FL 33487-2742

and by CRC Press
2 Park Square, Milton Park, Abingdon, OX14 4RN

© 2022 Taylor & Francis Group, LLC

CRC Press is an imprint of Taylor & Francis Group, LLC

Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, access www.copyright.com or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. For works that are not available on CCC please contact mpkbookspermissions@tandf.co.uk

Trademark notice: Product or corporate names may be trademarks or registered trademarks and are used only for identification and explanation without intent to infringe.

ISBN: 978-0-367-74845-6 (hbk)
ISBN: 978-1-003-15983-4 (ebk)

DOI: [10.1201/9781003159834](https://doi.org/10.1201/9781003159834)

Typeset in LM Roman
by KnowledgeWorks Global Ltd.

Contents

Preface	xv
1 Introduction to Statistical Science	1
1.1 Statistical Science: Description and Inference	1
1.1.1 Design, Descriptive Statistics, and Inferential Statistics	2
1.1.2 Populations and Samples	3
1.1.3 Parameters: Numerical Summaries of the Population	3
1.1.4 Defining Populations: Actual and Conceptual	4
1.2 Types of Data and Variables	4
1.2.1 Data Files	4
1.2.2 Example: The General Social Survey (GSS)	5
1.2.3 Variables	6
1.2.4 Quantitative Variables and Categorical Variables	6
1.2.5 Discrete Variables and Continuous Variables	7
1.2.6 Associations: Response Variables and Explanatory Variables	7
1.3 Data Collection and Randomization	8
1.3.1 Randomization	8
1.3.2 Collecting Data with a Sample Survey	9
1.3.3 Collecting Data with an Experiment	9
1.3.4 Collecting Data with an Observational Study	10
1.3.5 Establishing Cause and Effect: Observational versus Experimental Studies	10
1.4 Descriptive Statistics: Summarizing Data	11
1.4.1 Example: Carbon Dioxide Emissions in European Nations	11
1.4.2 Frequency Distribution and Histogram Graphic	11
1.4.3 Describing the Center of the Data: Mean and Median	13
1.4.4 Describing Data Variability: Standard Deviation and Variance	14
1.4.5 Describing Position: Percentiles, Quartiles, and Box Plots	15
1.5 Descriptive Statistics: Summarizing Multivariate Data	17
1.5.1 Bivariate Quantitative Data: The Scatterplot, Correlation, and Regression	17
1.5.2 Bivariate Categorical Data: Contingency Tables	18
1.5.3 Descriptive Statistics for Samples and for Populations	19
1.6 Chapter Summary	20
2 Probability Distributions	29
2.1 Introduction to Probability	29
2.1.1 Probabilities and Long-Run Relative Frequencies	29
2.1.2 Sample Spaces and Events	31
2.1.3 Probability Axioms and Implied Probability Rules	32
2.1.4 Example: Diagnostics for Disease Screening	33
2.1.5 Bayes' Theorem	34

2.1.6	Multiplicative Law of Probability and Independent Events	35
2.2	Random Variables and Probability Distributions	36
2.2.1	Probability Distributions for Discrete Random Variables	36
2.2.2	Example: Geometric Probability Distribution	37
2.2.3	Probability Distributions for Continuous Random Variables	38
2.2.4	Example: Uniform Distribution	38
2.2.5	Probability Functions (<i>pdf</i> , <i>pmf</i>) and Cumulative Distribution Function (<i>cdf</i>)	39
2.2.6	Example: Exponential Random Variable	40
2.2.7	Families of Probability Distributions Indexed by Parameters	41
2.3	Expectations of Random Variables	42
2.3.1	Expected Value and Variability of a Discrete Random Variable	42
2.3.2	Expected Values for Continuous Random Variables	43
2.3.3	Example: Mean and Variability for Uniform Random Variable	44
2.3.4	Higher Moments: Skewness	44
2.3.5	Expectations of Linear Functions of Random Variables	45
2.3.6	Standardizing a Random Variable	46
2.4	Discrete Probability Distributions	46
2.4.1	Binomial Distribution	46
2.4.2	Example: Hispanic Composition of Jury List	47
2.4.3	Mean, Variability, and Skewness of Binomial Distribution	48
2.4.4	Example: Predicting Results of a Sample Survey	49
2.4.5	The Sample Proportion as a Scaled Binomial Random Variable	50
2.4.6	Poisson Distribution	50
2.4.7	Poisson Variability and Overdispersion	51
2.5	Continuous Probability Distributions	52
2.5.1	The Normal Distribution	53
2.5.2	The Standard Normal Distribution	53
2.5.3	Examples: Finding Normal Probabilities and Percentiles	54
2.5.4	The Gamma Distribution	55
2.5.5	The Exponential Distribution and Poisson Processes	57
2.5.6	Quantiles of a Probability Distribution	57
2.5.7	Using the Uniform to Randomly Generate a Continuous Random Variable	58
2.6	Joint and Conditional Distributions and Independence	59
2.6.1	Joint and Marginal Probability Distributions	59
2.6.2	Example: Joint and Marginal Distributions of Happiness and Family Income	60
2.6.3	Conditional Probability Distributions	60
2.6.4	Trials with Multiple Categories: The Multinomial Distribution	61
2.6.5	Expectations of Sums of Random Variables	62
2.6.6	Independence of Random Variables	63
2.6.7	Markov Chain Dependence and Conditional Independence	64
2.7	Correlation between Random Variables	64
2.7.1	Covariance and Correlation	64
2.7.2	Example: Correlation between Income and Happiness	65
2.7.3	Independence Implies Zero Correlation, but Not Converse	66
2.7.4	Bivariate Normal Distribution *	66
2.8	Chapter Summary	69

3 Sampling Distributions	81
3.1 Sampling Distributions: Probability Distributions for Statistics	81
3.1.1 Example: Predicting an Election Result from an Exit Poll	81
3.1.2 Sampling Distribution: Variability of a Statistic's Value among Samples	83
3.1.3 Constructing a Sampling Distribution	84
3.1.4 Example: Simulating to Estimate Mean Restaurant Sales	85
3.2 Sampling Distributions of Sample Means	86
3.2.1 Mean and Variance of Sample Mean of Random Variables	86
3.2.2 Standard Error of a Statistic	87
3.2.3 Example: Standard Error of Sample Mean Sales	88
3.2.4 Example: Standard Error of Sample Proportion in Exit Poll	88
3.2.5 Law of Large Numbers: Sample Mean Converges to Population Mean	89
3.2.6 Normal, Binomial, and Poisson Sums of Random Variables Have the Same Distribution	89
3.3 Central Limit Theorem: Normal Sampling Distribution for Large Samples	90
3.3.1 Sampling Distribution of Sample Mean Is Approximately Normal	90
3.3.2 Simulations Illustrate Normal Sampling Distribution in CLT	92
3.3.3 Summary: Population, Sample Data, and Sampling Distributions	93
3.4 Large-Sample Normal Sampling Distributions for Many Statistics*	94
3.4.1 The Delta Method	95
3.4.2 Delta Method Applied to Root Poisson Stabilizes the Variance	96
3.4.3 Simulating Sampling Distributions of Other Statistics	96
3.4.4 The Key Role of Sampling Distributions in Statistical Inference	98
3.5 Chapter Summary	98
4 Statistical Inference: Estimation	105
4.1 Point Estimates and Confidence Intervals	105
4.1.1 Properties of Estimators: Unbiasedness, Consistency, Efficiency	106
4.1.2 Evaluating Properties of Estimators	107
4.1.3 Interval Estimation: Confidence Intervals for Parameters	107
4.2 The Likelihood Function and Maximum Likelihood Estimation	108
4.2.1 The Likelihood Function	108
4.2.2 Maximum Likelihood Method of Estimation	109
4.2.3 Properties of Maximum Likelihood (ML) Estimators	110
4.2.4 Example: Variance of ML Estimator of Binomial Parameter	111
4.2.5 Example: Variance of ML Estimator of Poisson Mean	111
4.2.6 Sufficiency and Invariance for ML Estimates	112
4.3 Constructing Confidence Intervals	113
4.3.1 Using a Pivotal Quantity to Induce a Confidence Interval	113
4.3.2 A Large-Sample Confidence Interval for the Mean	115
4.3.3 Confidence Intervals for Proportions	115
4.3.4 Example: Atheists and Agnostics in Europe	116
4.3.5 Using Simulation to Illustrate Long-Run Performance of CIs	117
4.3.6 Determining the Sample Size before Collecting the Data	117
4.3.7 Example: Sample Size for Evaluating an Advertising Strategy	118
4.4 Confidence Intervals for Means of Normal Populations	120
4.4.1 The t Distribution	120
4.4.2 Confidence Interval for a Mean Using the t Distribution	121
4.4.3 Example: Estimating Mean Weight Change for Anorexic Girls	122
4.4.4 Robustness for Violations of Normal Population Assumption	123

4.4.5	Construction of t Distribution Using Chi-Squared and Standard Normal	124
4.4.6	Why Does the Pivotal Quantity Have the t Distribution?	125
4.4.7	Cauchy Distribution: t Distribution with $df = 1$ Has Unusual Behavior	126
4.5	Comparing Two Population Means or Proportions	126
4.5.1	A Model for Comparing Means: Normality with Common Variability	127
4.5.2	A Standard Error and Confidence Interval for Comparing Means	127
4.5.3	Example: Comparing a Therapy to a Control Group	128
4.5.4	Confidence Interval Comparing Two Proportions	130
4.5.5	Example: Does Prayer Help Coronary Surgery Patients?	130
4.6	The Bootstrap	132
4.6.1	Computational Resampling and Bootstrap Confidence Intervals	132
4.6.2	Example: Booststrap Confidence Intervals for Library Data	132
4.7	The Bayesian Approach to Statistical Inference	134
4.7.1	Bayesian Prior and Posterior Distributions	135
4.7.2	Bayesian Binomial Inference: Beta Prior Distributions	136
4.7.3	Example: Belief in Hell	137
4.7.4	Interpretation: Bayesian versus Classical Intervals	138
4.7.5	Bayesian Posterior Interval Comparing Proportions	138
4.7.6	Highest Posterior Density (HPD) Posterior Intervals	138
4.8	Bayesian Inference for Means	139
4.8.1	Bayesian Inference for a Normal Mean	139
4.8.2	Example: Bayesian Analysis for Anorexia Therapy	140
4.8.3	Bayesian Inference for Normal Means with Improper Priors	141
4.8.4	Predicting a Future Observation: Bayesian Predictive Distribution .	142
4.8.5	The Bayesian Perspective, and Empirical Bayes and Hierarchical Bayes Extensions	142
4.9	Why Maximum Likelihood and Bayes Estimators Perform Well *	143
4.9.1	ML Estimators Have Large-Sample Normal Distributions	143
4.9.2	Asymptotic Efficiency of ML Estimators Same as Best Unbiased Estimators	145
4.9.3	Bayesian Estimators Also Have Good Large-Sample Performance .	146
4.9.4	The Likelihood Principle	146
4.10	Chapter Summary	147
5	Statistical Inference: Significance Testing	161
5.1	The Elements of a Significance Test	161
5.1.1	Example: Testing for Bias in Selecting Managers	161
5.1.2	Assumptions, Hypotheses, Test Statistic, P -Value, and Conclusion .	162
5.2	Significance Tests for Proportions and Means	164
5.2.1	The Elements of a Significance Test for a Proportion	164
5.2.2	Example: Climate Change a Major Threat?	166
5.2.3	One-Sided Significance Tests	166
5.2.4	The Elements of a Significance Test for a Mean	167
5.2.5	Example: Significance Test about Political Ideology	169
5.3	Significance Tests Comparing Means	170
5.3.1	Significance Tests for the Difference between Two Means	170
5.3.2	Example: Comparing a Therapy to a Control Group	171
5.3.3	Effect Size for Comparison of Two Means	172
5.3.4	Bayesian Inference for Comparing Two Means	173
5.3.5	Example: Bayesian Comparison of Therapy and Control Groups .	173

5.4	Significance Tests Comparing Proportions	174
5.4.1	Significance Test for the Difference between Two Proportions	174
5.4.2	Example: Comparing Prayer and Non-Prayer Surgery Patients	175
5.4.3	Bayesian Inference for Comparing Two Proportions	176
5.4.4	Chi-Squared Tests for Multiple Proportions in Contingency Tables . .	177
5.4.5	Example: Happiness and Marital Status	178
5.4.6	Standardized Residuals: Describing the Nature of an Association . .	179
5.5	Significance Test Decisions and Errors	180
5.5.1	The α -level: Making a Decision Based on the P -Value	181
5.5.2	Never “Accept H_0 ” in a Significance Test	181
5.5.3	Type I and Type II Errors	182
5.5.4	As $P(\text{Type I Error})$ Decreases, $P(\text{Type II Error})$ Increases	182
5.5.5	Example: Testing Whether Astrology Has Some Truth	184
5.5.6	The Power of a Test	185
5.5.7	Making Decisions versus Reporting the P -Value	186
5.6	Duality between Significance Tests and Confidence Intervals	186
5.6.1	Connection between Two-Sided Tests and Confidence Intervals	186
5.6.2	Effect of Sample Size: Statistical versus Practical Significance	187
5.6.3	Significance Tests Are Less Useful than Confidence Intervals	188
5.6.4	Significance Tests and P -Values Can Be Misleading	189
5.7	Likelihood-Ratio Tests and Confidence Intervals *	190
5.7.1	The Likelihood-Ratio and a Chi-Squared Test Statistic	191
5.7.2	Likelihood-Ratio Test and Confidence Interval for a Proportion . .	191
5.7.3	Likelihood-Ratio, Wald, Score Test Triad	192
5.8	Nonparametric Tests *	194
5.8.1	A Permutation Test to Compare Two Groups	194
5.8.2	Example: Petting versus Praise of Dogs	194
5.8.3	Wilcoxon Test: Comparing Mean Ranks for Two Groups	196
5.8.4	Comparing Survival Time Distributions with Censored Data	197
5.9	Chapter Summary	200
6	Linear Models and Least Squares	211
6.1	The Linear Regression Model and Its Least Squares Fit	211
6.1.1	The Linear Model Describes a Conditional Expectation	211
6.1.2	Describing Variation around the Conditional Expectation	212
6.1.3	Least Squares Model Fitting	213
6.1.4	Example: Linear Model for Scottish Hill Races	214
6.1.5	The Correlation	216
6.1.6	Regression toward the Mean in Linear Regression Models	217
6.1.7	Linear Models and Reality	218
6.2	Multiple Regression: Linear Models with Multiple Explanatory Variables .	219
6.2.1	Interpreting Effects in Multiple Regression Models	219
6.2.2	Example: Multiple Regression for Scottish Hill Races	220
6.2.3	Association and Causation	220
6.2.4	Confounding, Spuriousness, and Conditional Independence	221
6.2.5	Example: Modeling the Crime Rate in Florida	222
6.2.6	Equations for Least Squares Estimates in Multiple Regression	223
6.2.7	Interaction between Explanatory Variables in Their Effects	224
6.2.8	Cook’s Distance: Detecting Unusual and Influential Observations .	226
6.3	Summarizing Variability in Linear Regression Models	227
6.3.1	The Error Variance and Chi-Squared for Linear Models	228

6.3.2	Decomposing Variability into Model Explained and Unexplained Parts	228
6.3.3	<i>R</i> -Squared and the Multiple Correlation	229
6.3.4	Example: <i>R</i> -Squared for Modeling Scottish Hill Races	230
6.4	Statistical Inference for Normal Linear Models	231
6.4.1	The <i>F</i> Distribution: Testing That All Effects Equal 0	231
6.4.2	Example: Normal Linear Model for Mental Impairment	232
6.4.3	<i>t</i> Tests and Confidence Intervals for Individual Effects	233
6.4.4	Multicollinearity: Nearly Redundant Explanatory Variables	234
6.4.5	Confidence Interval for $E(Y)$ and Prediction Interval for Y	235
6.4.6	The <i>F</i> Test That All Effects Equal 0 is a Likelihood-Ratio Test *	236
6.5	Categorical Explanatory Variables in Linear Models	238
6.5.1	Indicator Variables for Categories	238
6.5.2	Example: Comparing Mean Incomes of Racial-Ethnic Groups	239
6.5.3	Analysis of Variance (ANOVA): An <i>F</i> Test Comparing Several Means	240
6.5.4	Multiple Comparisons of Means: Bonferroni and Tukey Methods	241
6.5.5	Models with Both Categorical and Quantitative Explanatory Variables	243
6.5.6	Comparing Two Nested Normal Linear Models	244
6.5.7	Interaction with Categorical and Quantitative Explanatory Variables	245
6.6	Bayesian Inference for Normal Linear Models	246
6.6.1	Prior and Posterior Distributions for Normal Linear Models	246
6.6.2	Example: Bayesian Linear Model for Mental Impairment	246
6.6.3	Bayesian Approach to the Normal One-Way Layout	247
6.7	Matrix Formulation of Linear Models *	248
6.7.1	The Model Matrix	248
6.7.2	Least Squares Estimates and Standard Errors	249
6.7.3	The Hat Matrix and the Leverage	250
6.7.4	Alternatives to Least Squares: Robust Regression and Regularization	250
6.7.5	Restricted Optimality of Least Squares: Gauss-Markov Theorem	250
6.7.6	Matrix Formulation of Bayesian Normal Linear Model	251
6.8	Chapter Summary	252
7	Generalized Linear Models	263
7.1	Introduction to Generalized Linear Models	263
7.1.1	The Three Components of a Generalized Linear Model	263
7.1.2	GLMs for Normal, Binomial, and Poisson Responses	264
7.1.3	Example: GLMs for House Selling Prices	265
7.1.4	The Deviance	267
7.1.5	Likelihood-Ratio Model Comparison Uses Deviance Difference	268
7.1.6	Model Selection: AIC and the Bias/Variance Tradeoff	268
7.1.7	Advantages of GLMs versus Transforming the Data	271
7.1.8	Example: Normal and Gamma GLMs for Covid-19 Data	271
7.2	Logistic Regression Model for Binary Data	272
7.2.1	Logistic Regression: Model Expressions	273
7.2.2	Parameter Interpretation: Effects on Probabilities and Odds	273
7.2.3	Example: Dose-Response Study for Flour Beetles	274
7.2.4	Grouped and Ungrouped Binary Data: Effects on Estimates and Deviance	276
7.2.5	Example: Modeling Italian Employment with Logit and Identity Links	278
7.2.6	Complete Separation and Infinite Logistic Parameter Estimates	279

7.3	Bayesian Inference for Generalized Linear Models	281
7.3.1	Normal Prior Distributions for GLM Parameters	281
7.3.2	Example: Logistic Regression for Endometrial Cancer Patients	282
7.4	Poisson Loglinear Models for Count Data	284
7.4.1	Poisson Loglinear Models	284
7.4.2	Example: Modeling Horseshoe Crab Satellite Counts	285
7.4.3	Modeling Rates: Including an Offset in the Model	286
7.4.4	Example: Lung Cancer Survival	287
7.5	Negative Binomial Models for Overdispersed Count Data *	288
7.5.1	Increased Variance Due to Heterogeneity	289
7.5.2	Negative Binomial: Gamma Mixture of Poisson Distributions	289
7.5.3	Example: Negative Binomial Modeling of Horseshoe Crab Data	290
7.6	Iterative GLM Model Fitting *	291
7.6.1	The Newton–Raphson Method	291
7.6.2	Newton–Raphson Fitting of Logistic Regression Model	292
7.6.3	Covariance Matrix of Parameter Estimators and Fisher Scoring	294
7.6.4	Likelihood Equations and Covariance Matrix for Poisson GLMs	294
7.7	Regularization with Large Numbers of Parameters *	295
7.7.1	Penalized Likelihood Methods	296
7.7.2	Penalized Likelihood Methods: The Lasso	297
7.7.3	Example: Predicting Opinions with Student Survey Data	297
7.7.4	Why Shrink ML Estimates toward 0?	299
7.7.5	Dimension Reduction: Principal Component Analysis	300
7.7.6	Bayesian Inference with a Large Number of Parameters	300
7.7.7	Huge n : Handling Big Data	300
7.8	Chapter Summary	301
8	Classification and Clustering	313
8.1	Classification: Linear Discriminant Analysis and Graphical Trees	314
8.1.1	Classification with Fisher’s Linear Discriminant Function	314
8.1.2	Example: Predicting Whether Horseshoe Crabs Have Satellites	314
8.1.3	Summarizing Predictive Power: Classification Tables and ROC Curves	316
8.1.4	Classification Trees: Graphical Prediction	318
8.1.5	Logistic Regression versus Linear Discriminant Analysis and Classification Trees	320
8.1.6	Other Methods for Classification: k -Nearest Neighbors and Neural Networks *	321
8.2	Cluster Analysis	324
8.2.1	Measuring Dissimilarity between Observations on Binary Responses .	325
8.2.2	Hierarchical Clustering Algorithm and Its Dendrogram	325
8.2.3	Example: Clustering States on Presidential Election Outcomes	326
8.3	Chapter Summary	328
9	Statistical Science: A Historical Overview	333
9.1	The Evolution of Statistical Science *	333
9.1.1	Evolution of Probability	333
9.1.2	Evolution of Descriptive and Inferential Statistics	334
9.2	Pillars of Statistical Wisdom and Practice	336
9.2.1	Stigler’s Seven Pillars of Statistical Wisdom	336
9.2.2	Seven Pillars of Wisdom for Practicing Data Science	338

Appendix A Using R in Statistical Science	341
A.0 Basics of R	341
A.0.1 Starting a Session, Entering Commands, and Quitting	341
A.0.2 Installing and Loading R Packages	341
A.0.3 R Functions and Data Structures	342
A.0.4 Data Input in R	344
A.0.5 R Control Flows	345
A.1 Chapter 1: R for Descriptive Statistics	345
A.1.1 Data Handling and Wrangling	345
A.1.2 Histograms and Other Graphics	346
A.1.3 Descriptive Statistics	347
A.1.4 Missing Values in Data Files	351
A.1.5 Summarizing Bivariate Quantitative Data	352
A.1.6 Summarizing Bivariate Categorical Data	353
A.2 Chapter 2: R for Probability Distributions	353
A.2.1 R Functions for Probability Distributions	353
A.2.2 Quantiles, Q-Q Plots, and the Normal Quantile Plot	355
A.2.3 Joint and Conditional Probability Distributions	358
A.3 Chapter 3: R for Sampling Distributions	358
A.3.1 Simulating the Sampling Distribution of a Statistic	358
A.3.2 Monte Carlo Simulation	359
A.4 Chapter 4: R for Estimation	361
A.4.1 Confidence Intervals for Proportions	361
A.4.2 Confidence Intervals for Means of Subgroups and Paired Differences .	362
A.4.3 The <i>t</i> and Other Probability Distributions for Statistical Inference .	362
A.4.4 Empirical Cumulative Distribution Function	363
A.4.5 Nonparametric and Parametric Bootstraps	364
A.4.6 Bayesian HPD Intervals Comparing Proportions	366
A.5 Chapter 5: R for Significance Testing	367
A.5.1 Bayes Factors and a Bayesian <i>t</i> Test	367
A.5.2 Simulating the Exact Distribution of the Likelihood-Ratio Statistic .	368
A.5.3 Nonparametric Statistics: Permutation Test and Wilcoxon Test . .	369
A.6 Chapter 6: R for Linear Models	370
A.6.1 Linear Models with the <code>lm</code> Function	370
A.6.2 Diagnostic Plots for Linear Models	370
A.6.3 Plots for Regression Bands and Posterior Distributions	371
A.7 Chapter 7: R for Generalized Linear Models	373
A.7.1 The <code>glm</code> Function	373
A.7.2 Plotting a Logistic Regression Model Fit	373
A.7.3 Model Selection for GLMs	373
A.7.4 Correlated Responses: Marginal, Random Effects, and Transitional Models	376
A.7.5 Modeling Time Series	377
A.8 Chapter 8: R for Classification and Clustering	379
A.8.1 Visualization of Linear Discriminant Analysis Results	379
A.8.2 Cross-Validation and Model Training	379
A.8.3 Classification and Regression Trees	381
A.8.4 Cluster Analysis with Quantitative Variables	381

Appendix B Using Python in Statistical Science	383
B.0 Basics of Python	383
B.0.1 Python Preliminaries	383
B.0.2 Data Structures and Data Input	384
B.1 Chapter 1: PYTHON for Descriptive Statistics	385
B.1.1 Random Number Generation	385
B.1.2 Summary Statistics and Graphs for Quantitative Variables	385
B.1.3 Descriptive Statistics for Bivariate Quantitative Data	386
B.1.4 Descriptive Statistics for Bivariate Categorical Data	388
B.1.5 Simulating Samples from a Bell-Shaped Population	388
B.2 Chapter 2: PYTHON for Probability Distributions	389
B.2.1 Simulating a Probability as a Long-Run Relative Frequency	389
B.2.2 Python Functions for Discrete Probability Distributions	390
B.2.3 Python Functions for Continuous Probability Distributions	391
B.2.4 Expectations of Random Variables	393
B.3 Chapter 3: PYTHON for Sampling Distributions	395
B.3.1 Simulation to Illustrate a Sampling Distribution	395
B.3.2 Law of Large Numbers	395
B.4 Chapter 4: PYTHON for Estimation	396
B.4.1 Confidence Intervals for Proportions	396
B.4.2 The <i>t</i> Distribution	396
B.4.3 Confidence Intervals for Means	396
B.4.4 Confidence Intervals Comparing Means and Comparing Proportions .	397
B.4.5 Bootstrap Confidence Intervals	398
B.4.6 Bayesian Posterior Intervals for Proportions and Means	399
B.5 Chapter 5: PYTHON for Significance Testing	400
B.5.1 Significance Tests for Proportions	400
B.5.2 Chi-Squared Tests Comparing Multiple Proportions in Contingency Tables	400
B.5.3 Significance Tests for Means	401
B.5.4 Significance Tests Comparing Means	401
B.5.5 The Power of a Significance Test	403
B.5.6 Nonparametric Statistics: Permutation Test and Wilcoxon Test . . .	403
B.5.7 Kaplan-Meier Estimation of Survival Functions	404
B.6 Chapter 6: PYTHON for Linear Models	404
B.6.1 Fitting Linear Models	404
B.6.2 The Correlation and R-Squared	406
B.6.3 Diagnostics: Residuals and Cook's Distances for Linear Models . .	407
B.6.4 Statistical Inference and Prediction for Linear Models	410
B.6.5 Categorical Explanatory Variables in Linear Models	411
B.6.6 Bayesian Fitting of Linear Models	412
B.7 Chapter 7: PYTHON for Generalized Linear Models	413
B.7.1 GLMs with Identity Link	413
B.7.2 Logistic Regression: Logit Link with Binary Data	415
B.7.3 Separation and Bayesian Fitting in Logistic Regression	416
B.7.4 Poisson Loglinear Model for Counts	417
B.7.5 Negative Binomial Modeling of Count Data	420
B.7.6 Regularization: Penalized Logistic Regression Using the Lasso . .	421
B.8 Chapter 8: PYTHON for Classification and Clustering	421
B.8.1 Linear Discriminant Analysis	421
B.8.2 Classification Trees and Neural Networks for Prediction	423

B.8.3 Cluster Analysis	425
Appendix C Brief Solutions to Exercises	427
C.1 Chapter 1: Solutions to Exercises	427
C.2 Chapter 2: Solutions to Exercises	429
C.3 Chapter 3: Solutions to Exercises	431
C.4 Chapter 4: Solutions to Exercises	433
C.5 Chapter 5: Solutions to Exercises	436
C.6 Chapter 6: Solutions to Exercises	439
C.7 Chapter 7: Solutions to Exercises	443
C.8 Chapter 8: Solutions to Exercises	446
Bibliography	447
Example Index	449
Subject Index	453

Preface

This book presents an overview of the *foundations*—the key concepts and results—of statistical science. The primary intended audience is undergraduate students who are training to become data scientists. This is not a book, however, about how to become a data scientist or about the newest methods being used by data scientists or about how to analyze “big data” and the wide variety of types of data with which data scientists deal. It is a book that has the purpose of teaching potential data scientists the foundations of one of the core tenets of data science—statistical science.

Statistical science is by now a large subject, having many distinct specialties. This book highlights the topics with which we believe that any data scientist should be familiar: descriptive statistical methods, probability distributions, the inferential statistical methods of confidence intervals and significance testing, and linear and generalized linear modeling. When combined with courses that a student majoring in Data Science would take in computer science and mathematics, this book provides the background needed to follow this introduction to statistical science by studying specialized areas of it.¹

This book assumes that students have knowledge of calculus, so we can focus on *why a statistical analysis works* as well as *how to do it*. University statistical science courses that require a calculus background often have the name *Mathematical Statistics*. We avoid this term in the title of our book, as we do not want students to think that statistical science is a subfield of mathematics or that complex mathematics is necessary to be proficient in understanding and applying statistical science. In fact, we mainly use only basic calculus tools of differentiation and integration, and then only for some topics. Compared to the content of traditional mathematical statistics textbooks, our book has less emphasis on probability theory, derivations of probability distributions of transformations of random variables, decision theory, and statement and formal proof of theorems. It introduces some modern topics that do not normally appear in such texts but are especially relevant for data scientists, such as generalized linear models for non-normal responses, Bayesian and regularized model-fitting, and classification and clustering. The greatest difference from a traditional mathematical statistics book, however, is that this book shows how to implement statistical methods with modern software and illustrates statistical concepts and theory using simulations.

To use and properly interpret methods of modern statistical science, computational skills are as important as mathematical skills. Besides using mathematics to show “why it works,” we use computational simulations and Internet apps to help provide intuition about foundational results such as behavior of sampling distributions and error rates for statistical inferences. Throughout the book, examples with real data show how to use the free statistical software **R** to implement statistical methods. The book also contains software appendices that present greater detail about **R** as well as introducing **Python** for statistical analyses. The **Python** appendix shows analyses for the examples analyzed in the chapters with **R**, so an instructor can easily use the text in a course that has **Python** as its primary software. Since

¹Such as multivariate analysis, nonparametrics, categorical data analysis, design and sample survey methods, time series, longitudinal data analysis, survival analysis, decision theory, Bayesian statistics, stochastic modeling, computational methods of statistics, and smoothing and nonlinear modeling

the book focuses on the foundations of statistical science, it has less emphasis on some practical issues of data analysis, such as preparing and cleaning data files. However, the software appendices also introduce additional analyses that supplement the examples presented in the chapters. A regularly-updated website <http://stat4ds.rwth-aachen.de> for the book has all data files analyzed and longer versions of the software appendices as well as an appendix about the use of **Matlab** for statistical analyses. The data files are also available at www.stat.ufl.edu/~aa/ds and at the GitHub site <https://github.com/stat4DS/data>.

Use of this book as a course textbook

Chapters 1–6 of this book are designed as a textbook for an introductory course on statistical science for undergraduate students majoring in Data Science or Statistics or Mathematics. Some instructors may prefer to skip some of the less central or more technical material, such as Sections 3.4, 4.9, 5.7, 5.8, and 6.7. (These and other such sections and subsections have an * next to their titles.) With all nine chapters and the extra material presented in the R and Python appendices, it is also appropriate for a two-term sequence of courses. The book also can serve programs that have a heavy focus on statistical science, such as econometrics and operations research. It also should be useful to graduate students in the social, biological, and environmental sciences who choose Statistics as their minor area of concentration, so they can learn about the foundations that underlie statistical methods that they use. An instructor can use either R or Python as the main software for the course, as the examples in the main part of the text use R but the same examples are shown with Python in its appendix.

Each chapter contains many exercises for students to practice and extend the theory and methods. The exercises are grouped into two parts: Exercises in *Data Analysis and Applications* request that students perform data analyses similar to the ones presented in that chapter. Exercises in *Methods and Concepts* relate directly to the foundations aspect of the book. They ask questions about properties of statistical methods, conceptual questions about their bases, as well as extend that chapter’s results. An appendix contains outlines of solutions for the odd-numbered exercises.

This book is by no means a complete overview of statistical science. The field is large and grows more every year, with areas being developed now that did not even exist in the twentieth century. However, we do believe it provides a solid introduction to the core material with which we believe any data scientist should be familiar.

In preparing this book together, Agresti (agresti@ufl.edu) has taken main responsibility for the chapter material and Kateri (maria.kateri@rwth-aachen.de) has taken main responsibility for the appendices about R and Python statistical software and the expanded appendices about R, Python, and Matlab at the book’s website. We welcome any comments or suggestions that you care to send either of us that we can take into account in future editions of this book.

Acknowledgments

Thanks to several friends and colleagues who provided comments on various versions of this manuscript or who provided data sets or other help, including Alessandra Brazzale, Jane Brockmann, Brian Caffo, Sir David Cox, Bianca De Stavola, Cristina Cuesta, Travis Gerke, Sabrina Giordano, Anna Gottard, Ralitza Gueorguieva, Bernhard Klingenberg, Bhramar Mukherjee, Ranjini Natarajan, Madan Oli, Euijung Ryu, Alessandra Salvan, Nicola Sartori, Elena Stanghellini, Stephen Stigler, Gerhard Tutz, Roberta Varriale, Larry Winner, and Daniela Witten. Thanks to Hassan Satvat for help in setting up the book’s website and to Bernhard Klingenberg for developing the excellent apps at www.artofstat.com that are often cited in the book. Many thanks to Joyce Robbins, Mintaek Lee, Jason M. Graham, Christopher Gaffney, Tumulesh Solanky, and Steve Chung for providing helpful reviews to

CRC Press of our manuscript. Finally, special thanks to John Kimmel, Executive Editor of Statistics for Chapman & Hall/CRC Press, for his encouragement and support in this book project.

ALAN AGRESTI and MARIA KATERI
*Gainesville Florida and Brookline Massachusetts, USA;
Aachen, Germany*

April 2021



Taylor & Francis
Taylor & Francis Group
<http://taylorandfrancis.com>

1

Introduction to Statistical Science

Compared to mathematics and the physical and natural sciences, statistical science is quite young. The statistical methods that you'll learn about in this book were mainly developed within the past century. Modern computing power is causing a revolution in the sorts of data analyses that are possible, so new methods are continually being developed. In recent years, new statistical methods have resulted from challenges in analyzing data in diverse fields, such as medicine (e.g., genetic data relating to disease presence, data for personalized medical decisions) and business (e.g., data on consumer buying behavior, data from experiments comparing advertising strategies). This book presents the foundations underlying the methods of statistical science, explaining when and why these methods work, and shows how to use statistical software to apply them.

Statistical software also has become increasingly powerful and easily available. This has had benefits in the data analyses that are now possible, but a danger is that prospective data scientists might think that statistical science is merely a computational toolbox consisting of a variety of algorithms. A goal of this book, by contrast, is to show that *the methods of statistical science result from a unified theory*, although that theory itself has slight variations in the way it is implemented or interpreted. Another danger of the ubiquity of statistical software is that prospective data scientists might expect that software can automatically perform good data analyses without input from the user. We'll see, however, that careful thought is needed to decide which statistical methods are appropriate for any particular situation, as they all make certain assumptions, and some methods work poorly when the assumptions are violated. Moreover, a data scientist needs to be able to interpret and explain the results that software yields.

In this chapter, we introduce statistical science as a field that deals with describing data and using them to make inferences. We define types of *variables* that represent how measured characteristics can vary from observation to observation. We also introduce *graphical and numerical methods* for describing the data. When a study can use *randomization* in collecting the data or conducting an experiment, data analysts can exploit the random variation to make reliable estimations and predictions.

1.1 Statistical Science: Description and Inference

You already have a sense of what the word *statistics* means. You regularly hear statistics quoted about sports events, the economy, medical research, and opinions, beliefs, and behaviors of people. In this sense, a statistic is merely a number calculated from *data* —the observations that provide information about the subject matter. But the field of *statistical science* has a much broader sense—as a field that gives us a way of gathering and analyzing the data in an objective manner.

Statistical science

Statistical science is the science of developing and applying methods for collecting, analyzing, and interpreting data.

Many methods of statistical science incorporate reasoning using tools of *probability*. The methods enable us to deal with uncertainty and variability in virtually all scientific fields. With statistical methods, we learn from the data while measuring, controlling, and communicating uncertainty.

1.1.1 Design, Descriptive Statistics, and Inferential Statistics

Statistical science has three aspects:

1. **Design:** Planning how to gather relevant data for the subject matter of interest.
2. **Description:** Summarizing the data.
3. **Inference:** Making evaluations, such as estimations and predictions, based on the data.

Design refers to planning a study so that it yields useful data. For example, for a poll taken to determine public opinion on some issue, the design specifies how to select the people to interview and constructs the questionnaire for interviews. For a research study to compare an experimental diet with a standard diet to address obesity, the design specifies how to obtain people for the study, how to determine which people use each diet, and specifies the characteristics to measure to compare the diets.

Description refers to summarizing the data, to mine the information that the data provide. For any study, the raw data are a complete listing of observations that can be overwhelming for comprehension. To present the results, we reduce the data to simpler and more understandable form without distorting or losing much information. Graphs, tables, and numerical summaries such as averages and percentages are called **descriptive statistics**.

Inference refers to using the data to make estimations and other sorts of evaluations, such as predictions. These evaluations take into account random variability that occurs with the characteristics measured and the resulting uncertainty in decision-making. For instance, suppose that in the study comparing two diets, the people on the experimental diet had an average weight loss of 7.0 kilograms. What can we say about the average weight change if hypothetically *all* obese people used this diet? An inferential statistical method provides an interval of numbers within which we can predict that the average weight change would fall. The analysis might enable us to conclude, with small probability of being wrong, that the average weight change for all obese people would fall between 5.6 and 8.4 kilograms. Another inferential statistical method would enable us to decide whether that average weight change is greater than would be obtained with a standard diet or no special diet. Other inferential statistical methods evaluate whether weight change is associated with characteristics other than the diet, such as a person's gender, race, age, attained education, and amount of weekly exercise. Data-based evaluations such as estimations and predictions are called **statistical inferences**.

Descriptive statistics and *inferential statistics* are the two main types of methods for analyzing data. Researchers use them to answer questions such as, “Does the experimental diet have a beneficial effect in reducing obesity, and is it more effective than a standard diet?” “How does the sales of a product compare if we place an advertisement for it at websites, or in mailings, or in newspapers, or on TV programs?” “Do states in the U.S.

that have stronger gun control laws tend to have lower murder rates, taking into account socioeconomic factors?” “Is student performance in Canada associated with the amount of money spent per student, the size of the classes, or the teachers’ salaries?” “Do a majority of all New Zealanders favor legalization of marijuana?”

1.1.2 Populations and Samples

The entities on which a study makes observations are called the sample *subjects*. Usually the subjects are individual people, such as in a survey, but they need not be. For example, an agricultural experiment might have cows as subjects if its goal is to compare milk yields for different diets. An ecological survey might have different forest areas as subjects in a study of species diversity. Subjects in social surveys might be people, families, schools, or counties.

Although we obtain data for the sample subjects, our ultimate interest is on the *population* from which the sample is taken.

Population and sample

The **population** is the total set of subjects of interest. A **sample** is the set of subjects from the population for which data are available.

The goal of most data analyses is to learn about populations. But it is almost always necessary, and more practical, to observe only samples from those populations. For example, polling organizations such as the Gallup poll (www.gallup.com) and the Pew Research Center (www.pewresearch.org) usually sample about 1000–2000 Americans to gather information about opinions and beliefs of the population of *all* adult Americans.

Inferential statistics provide evaluations about a population, based on data from a sample. For example, a survey taken in the U.S. in 2018 asked, “Do you believe in heaven?” The population of interest was all adults in the United States. Of the 1141 sampled subjects, 81% answered *yes*. We would be interested, however, not only in those 1141 people but in the *population* of more than 250 million adults in the U.S. An inferential method presented in Chapter 4 estimates that the population percentage believing in heaven almost certainly falls between 78% and 84%. That is, the sample value of 81% has a “margin of error” of 3%. Inferential statistical analyses can predict characteristics of entire populations quite well by selecting samples that are very small relative to the population size. In this book, we’ll learn why this works.

1.1.3 Parameters: Numerical Summaries of the Population

To distinguish between a descriptive statistic calculated for a sample and the corresponding characteristic of the population, we use the term **parameter** for the population characteristic.

Parameter

A **parameter** is a numerical summary of a population.

In practice, our primary interest is in the values of parameters rather than sample descriptive statistics. For example, in viewing the results of a poll before an election, we would be more interested in the *population* percentages favoring the various candidates than in the *sample* percentages for the people interviewed. However, parameter values are almost

always unknown. The sample and statistics describing it help us to make inferences about the unknown parameter values.

A key aspect of statistical inference involves reporting the *precision* of the sample statistic that estimates the population parameter. For the example on belief in heaven, the reported margin of error of 3% predicted how close the *sample* value of 81% was to the unknown *population* percentage. The other key aspect relates to the *probability* with which we obtain that precision. For instance, a statistical inference might state that we can be 95% sure that the sample value of 81% differs from the population value by no more than 3%.

1.1.4 Defining Populations: Actual and Conceptual

Usually the population to which inferences apply is an actual set of subjects, such as all adult residents of a nation. Sometimes, though, the inferences refer to a *conceptual* population—one that does not actually exist but is hypothetical.

For example, suppose a medical research team investigates a newly proposed drug for treating a virus by conducting a study at several medical centers. Such a medical study is called a *clinical trial*. The study would compare virus patients who are given the new drug to other virus patients who instead receive a standard treatment or a placebo, using descriptive statistics such as the percentages who respond positively. In applying inferential statistical methods, the researchers would like their inferences to apply to the conceptual population of *all* people suffering from the virus now or at some time in the future.

1.2 Types of Data and Variables

The observations gathered on the characteristics of interest are the *data*. For example, a survey of 1000 people to analyze opinions about the legalization of same-sex marriage might also observe characteristics such as political party affiliation, frequency of attending religious services, number of years of education, annual income, marital status, race, and gender. The data for a particular subject would consist of observations such as (opinion = do not favor legalization, political party = Republican, religiosity = attend services once a week, education = 12 years, annual income in the interval 40–60 thousand dollars, marital status = married, race = White, gender = male).

1.2.1 Data Files

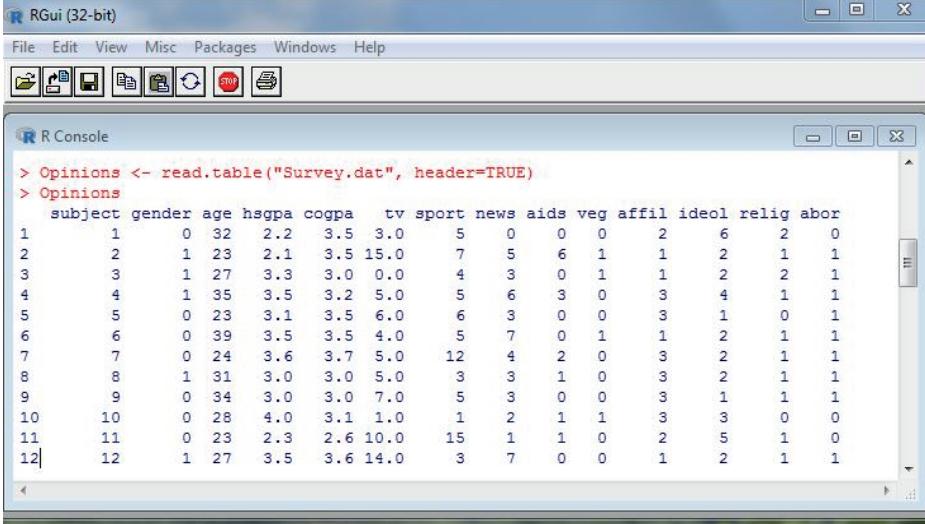
Statistical software analyzes data organized in the spreadsheet form of a *data file*:

- Any one row of a data file contains the observations for a particular subject (e.g., person) in the sample.
- Any one column of a data file contains the observations for a particular characteristic (e.g., opinion about legalized same-sex marriage).

The *number* of subjects in the data file, called the *sample size*, is denoted by n . If we observe 10 characteristics for a sample of $n = 2000$ people, then the data file has 2000 rows and 10 columns.

Throughout this book, we use the statistical software R to illustrate statistical analyses. This software package is available to download for free at www.r-project.org. In R, we

save a data file as a *data frame*,¹ which is the fundamental data structure required by many R functions. [Figure 1.1](#) shows the use of R to read a data file called **Survey.dat** (containing data from a student survey mentioned in Exercise 1.2) from a directory on a computer, save it as a data frame called *Opinions*, and display it.



```
> Opinions <- read.table("Survey.dat", header=TRUE)
> Opinions
  subject gender age hsgpa cogpa   tv sport news aids veg affil ideol relig abor
  1       1     0  32   2.2   3.5  3.0    5   0   0   0   2     6    2    0
  2       2     1  23   2.1   3.5 15.0    7   5   6   1   1     2    1    1
  3       3     1  27   3.3   3.0  0.0    4   3   0   1   1     2    2    1
  4       4     1  35   3.5   3.2  5.0    5   6   3   0   3     4    1    1
  5       5     0  23   3.1   3.5  6.0    6   3   0   0   3     1    0    1
  6       6     0  39   3.5   3.5  4.0    5   7   0   1   1     2    1    1
  7       7     0  24   3.6   3.7  5.0   12   4   2   0   3     2    1    1
  8       8     1  31   3.0   3.0  5.0    3   3   1   0   3     2    1    1
  9       9     0  34   3.0   3.0  7.0    5   3   0   0   3     1    1    1
 10      10    0  28   4.0   3.1  1.0    1   2   1   1   3     3    0    0
 11      11    0  23   2.3   2.6 10.0   15   1   1   0   2     5    1    0
 12      12    1  27   3.5   3.6 14.0    3   7   0   0   1     2    1    1
```

FIGURE 1.1 Part of a R session for loading and displaying a data file.

Existing archived collections of data are called **databases**. Many databases result from a survey or study of some type, but some are existing records of data that result from other purposes. An example is a database of patients' electronic medical records. With the increasing variety of data that can be recorded electronically, not all data files have the format of a traditional data file with entries that are numbers or characters. For example, in medical records, some observations for some subjects may be images, such as a mammogram, a chest x-ray, or a brain scan,² or a continuous streaming of data over time, such as monitoring of heart-rate, respiratory-rate, blood pressure, and temperature.

1.2.2 Example: The General Social Survey (GSS)

Some databases are freely available on the Internet. An important database in the U.S. contains results since 1972 of the *General Social Survey* (GSS), conducted every other year by the National Opinion Research Center at the University of Chicago. It gathers information using personal interviews of a sample of about $n = 2000$ subjects from the U.S. adult population to provide a snapshot of opinions and behaviors. Researchers use it to investigate how adult Americans answer a wide diversity of questions, such as, “Do you believe in life after death?” and “Would you be willing to pay higher prices in order to protect the environment?” Similar social surveys occur in other countries, such as the General Social Survey administered by Statistics Canada, the British Social Attitudes Survey, and the Eurobarometer survey and European Social Survey for nations in the European Union.

It is easy to get summaries of data from the GSS database:

¹For details, see [Section A.0.3](#) of the R Appendix.

²See <https://aimi.stanford.edu/research/public-datasets> for examples of data files of this type.

- Go to the website <https://sda.berkeley.edu/archive.htm> at the Survey Documentation and Analysis site at the University of California, Berkeley.
- Click on the most recently available *General Social Survey (GSS) Cumulative Datafile*. You will then see a “variable selection” listing in the left margin of characteristics measured over the years, and a menu on the right for selecting particular characteristics of interest.
- Type the name of a characteristic of interest in the *Row* box, and click on *Run the table*. The GSS site will then generate a table that shows the possible values for the characteristic and the number of people and the percentage who made each possible response.

For example, in one survey the GSS asked “About how many good friends do you have?” The GSS name for this characteristic is NUMFREND. The table that the GSS provides shows that the responses of 1, 2, 3, 4, 5, and 6 good friends had the percentages 6.1, 16.2, 15.7, 14.2, 11.3, and 8.8, respectively, with the remaining 27.7% spread around the other possible responses.

1.2.3 Variables

For the characteristics we measure in a study, *variability* occurs naturally among subjects in a sample or population. For instance, variation occurs from student to student in their college grade point average (GPA). A study to investigate the factors mainly responsible for that variability might also observe other characteristics that vary among students, such as high school GPA, college board score, time per day spent studying, time per day watching TV or browsing the Internet, and whether at least one parent attended college. Any characteristic that we can measure for the subjects is called a **variable**. The term reflects that values of the characteristic *vary* among subjects.

Variable

A **variable** is a characteristic that can vary in value among subjects in a sample or population.

The values the variable can take form a *measurement scale*. The valid statistical methods for a variable depend on its measurement scale. We treat a numerical-valued variable such as number of good friends differently than a variable measured with categories, such as (*yes*, *no*) for whether employed. We next present two ways to classify variables. The first type refers to whether the measurement scale consists of numbers or categories. The second type refers to the fineness of measurement—the number of values in the measurement scale.

1.2.4 Quantitative Variables and Categorical Variables

A variable is called **quantitative** when the measurement scale has numerical values that represent different magnitudes of the variable. Examples of quantitative variables are number of good friends, annual income, college GPA, age, and weight.

A variable is called **categorical** when the measurement scale is a set of categories. Examples of categorical variables are marital status (with categories such as *single*, *married*, *divorced*, *widowed*), primary mode of transportation to work (*automobile*, *bicycle*, *bus*, *subway*, *walk*), preferred destination for clothes shopping (*downtown*, *Internet*, *mall*, *other*), and favorite type of music (*classical*, *country*, *folk*, *jazz*, *rap/hip-hop*, *rock*). Categorical variables having only two categories, such as whether employed (*yes*, *no*), are called **binary**.

For categorical variables, distinct categories differ in quality, not in numerical magnitude. Categorical variables are often called ***qualitative***.

Categorical variables have two types of measurement scales. For some categorical variables, such as the ones just mentioned, the categories are unordered. The scale does not have a “high” or “low” end. The categories are then said to form a ***nominal scale***. By contrast, some categorical scales have a natural *ordering* of values. The categories form an ***ordinal scale***. Examples are perceived happiness (*not too happy, pretty happy, very happy*), headache pain (*none, slight, moderate, severe*), and political philosophy (*very liberal, slightly liberal, moderate, slightly conservative, very conservative*).

1.2.5 Discrete Variables and Continuous Variables

Another classification refers to the *number* of values in the measurement scale.

Discrete and continuous variables

A quantitative variable is ***discrete*** if it can take a set of distinct, separate values, such as the nonnegative integers (0, 1, 2, 3, ...). It is ***continuous*** if it can take an infinite continuum of possible real number values.

Examples of discrete variables are one’s number of good friends, number of computers in household, and number of days playing a sport in the past week. Any variable phrased as “the number of …” is discrete, because we can list its possible values (0, 1, 2, 3, ...). Examples of continuous variables are height, weight, age, distance a person walks in a day, winning time in a marathon race, and how long a cell phone works before it needs recharging. It is impossible to write down all the distinct potential values, because they form an interval of infinitely many real-number values. A person’s age, for example, could take the value 20.6294473… years.

In practice, we round continuous variables when measuring them, so the actual measurement is discrete. We say that an individual is 20 years old whenever that person’s age is somewhere between 20 and 21. On the other hand, some variables, although discrete, have a very large number of possible values. In measuring annual income in dollars, the potential values are 0, 1, 2, 3, ..., up to some very large value in many millions. Statistical methods for continuous variables are used for quantitative variables that can take a very large number of values, regardless of whether they are theoretically continuous or discrete.

1.2.6 Associations: Response Variables and Explanatory Variables

Most studies have more than one variable. With multivariable analyses, we say that an ***association*** occurs between two variables if certain values of one variable tend to go with certain values of the other. For example, consider religious affiliation, with categories (*Catholic, Protestant, Muslim, Jewish, Other*), and ethnic group, with categories (*African-American, Anglo-American, Hispanic, Other*). In the United States, Anglo-Americans are more likely to be Protestant than are Hispanics, who are overwhelmingly Catholic. African-Americans are even more likely to be Protestant. An association exists between religious affiliation and ethnic group, because the percentage of people having a particular religious affiliation changes as ethnic group changes.

When we study the association between two variables, usually one is an outcome variable on which comparisons are made at levels of the other variable. The outcome variable is called the ***response variable***. The variable that defines the groups is called the ***explanatory variable***. The analysis studies how the outcome on the response variable *depends on* or

is *explained by* the value of the explanatory variable. For example, when we describe how religious affiliation depends on ethnic group, religious affiliation is the response variable. In a comparison of men and women on annual income, annual income is the response variable and gender is the explanatory variable.

Most studies have one response variable and multiple explanatory variables. Sometimes the response variable is called the ***dependent variable*** and the explanatory variables are called the ***independent variables***. We prefer not to use these terms, because the terms *independent* and *dependent* are used in statistical science for many other things, and the terminology suggests causal interpretations that are usually inappropriate.

1.3 Data Collection and Randomization

When we apply inferential methods of statistical science to parameters for some population, the quality of the inferences depends on how well the *sample* represents the *population*.

1.3.1 Randomization

Randomization is a mechanism for achieving good sample representation of a population in a survey or an experiment.

Simple random sample

A ***simple random sample*** of n subjects from a population is one in which each possible sample of size n has the same probability (chance) of being selected.

With simple random sampling, everyone has the same chance of inclusion in the sample, so it is fair. It tends to yield a sample that resembles the population. This reduces the chance that the sample is seriously biased in some way, leading to inaccurate inferences about the population.

Suppose that a researcher in a medical center plans to compare two drugs for some adverse condition. She has four patients with this condition, and she wants to randomly select two to use each drug. Denote the four patients by P_1 , P_2 , P_3 , and P_4 . In selecting $n = 2$ subjects to use the first drug, the six possible samples are

$$(P_1, P_2), (P_1, P_3), (P_1, P_4), (P_2, P_3), (P_2, P_4), (P_3, P_4).$$

More generally, let N denote the population size. The population has $\binom{N}{n}$ possible samples of size n . For example, a population of size $N = 4$ has $\binom{4}{2} = 4!/[2!(4-2)!] = 6$ possible samples of size $n = 2$. You could select the simple random sample by placing the four people's names on four identical ballots and selecting two blindly from a hat. This is unwieldy with the larger values for N and n usual in practice, and these days software can easily select the sample from a list of the population members using a *random number generator*.

We illustrate for random selection of $n = 5$ students out of a class of size $N = 60$. We assign the numbers $01, 02, \dots, 60$ to the class members and generate five random numbers between 01 and 60. With the statistical software R, the `sample` function performs simple random sampling from a numbered population list:³

³We suggest that you read the Basics and Chapter 1 sections in this book's R Appendix to learn more about R and its use for descriptive statistical analysis.

```
| > sample(1:60, 5)    # Comments about R commands follow the # symbol  
| [1] 11 55 48 59 29  # output line [1] shows the five integers randomly generated
```

The sample of size 5 selects the students numbered 11, 55, 48, 59, 29.

The *simple* adjective in “simple random sample” distinguishes this type of sampling from more complex sampling schemes that also have elements of randomization. For instance, *stratified random sampling* schemes divide the population into separate groups (“strata”) and then select a simple random sample from each stratum. This is useful for comparing groups on some variable when a particular group is relatively small and may not be adequately represented in a simple random sample. *Cluster random sampling* schemes divide the population into a large number of clusters, such as city blocks, and select a simple random sample of the clusters. This is useful when a complete listing of the population is not available.

1.3.2 Collecting Data with a Sample Survey

Some studies sample people from a population and interview them to collect data. This method of data collection is called a *sample survey*. The interview could be a personal interview, telephone interview, or self-administered questionnaire. Implementing a sample survey requires finding or constructing a valid *population list* for selecting the sample. A major challenge with collecting data with sample surveys is *nonresponse*—a significant percentage of those sampled refuse to be interviewed or don’t respond on some items.

The *General Social Survey*, introduced in [Section 1.2.2](#), is a sample survey. The GSS uses a random sampling design that is more complex than simple random sampling, incorporating multiple stages and clustering to make the survey easier to implement, but it ensures adequate coverage by giving each family the same chance of inclusion.

1.3.3 Collecting Data with an Experiment

Some studies use a planned *experiment* to generate data. An experiment compares subjects on a response variable under different conditions. Those conditions, which are levels of an explanatory variable, are called *treatments*. For instance, the treatments might be different drugs for treating some illness, compared in a clinical trial. The researcher specifies a plan for how to assign subjects to the treatments, called the *experimental design*. Good experimental designs use randomization to determine which treatment a subject receives. This reduces bias and allows us to use statistical inference.

For example, the Physicians’ Health Study Research Group at Harvard Medical School designed an experiment to analyze whether regular intake of aspirin reduces mortality from heart disease. Of about 22,000 physicians, half were randomly chosen to take an aspirin every other day. The remaining half took a placebo, which had no active agent. After five years, rates of heart attack were compared. By using randomization to determine who received which treatment, the researchers knew the groups would roughly balance on all variables that could affect heart attack rates, such as age and quality of health. If the physicians could decide on their own which treatment to take, the groups might have been out of balance on some important factor. For instance, if younger physicians were more likely to select aspirin, then a lower heart attack rate among the aspirin group could occur merely because younger subjects are less likely to suffer heart attacks.

In medical research, *randomized clinical trials* are experiments using randomization that have been the gold standard for many years. But experiments are now used to address questions in an increasing variety of areas. For instance, the economists at Harvard and Massachusetts Institute of Technology (MIT) who won the Nobel Prize in 2019 pioneered

the use of experiments to determine the policies that best improve the lives of the poor.⁴ They randomly selected participants for anti-poverty programs, such as to evaluate if access to textbooks or access to remedial instruction improves education results. Randomized experiments have addressed diverse topics such as how to improve child nutrition, protect forests, and reduce gender discrimination.

1.3.4 Collecting Data with an Observational Study

In many application areas, it is not possible to conduct experiments to answer the questions of interest. We cannot randomly assign subjects to the groups we want to compare, such as levels of gender or race or educational level or annual income or usage of guns. Many studies merely *observe* the outcomes for available subjects on the variables of interest, without any experimental control of the subjects. Such studies are called ***observational studies***. Sample surveys are examples of observational studies.

Sometimes we can envision an experiment that would help us answer some question, but the experiment would be *unethical* to conduct, so observational studies are used instead. For example, in the mid-20th century some researchers decided to investigate whether an association exists between lung cancer and smoking. The researchers could have answered the question by taking a group of youngsters, randomly splitting them into two groups, and instructing one group to smoke a pack of cigarettes each day and the other group not to smoke at all. Then, after 50 years, the study would have compared the percentages of smokers and non-smokers who got lung cancer. Such an experiment would not have been ethical or feasible to conduct, and the answer was needed at that time itself and not 50 years later, so it was necessary to use observational studies to address this issue. One useful approach sampled hospital patients, matching each adult suffering from lung cancer to an adult control of similar age who did not have it,⁵ and compared the lung-cancer group with the control group in terms of how much they had smoked in the past.

1.3.5 Establishing Cause and Effect: Observational versus Experimental Studies

With observational studies, making causal conclusions based on comparing groups on a response variable is dangerous because the groups may be imbalanced on other variables that affect the response outcome. This is true even with random sampling. For instance, suppose we plan to compare performance on some standardized exam for Black, Hispanic, and White students. If White students have a higher average score, a variety of variables might account for that difference. Perhaps, on the average, White students have higher parents' attained education or higher parents' income or better quality of school attended. Those or other key variables may not even have been measured in the study.

Establishing *cause and effect* definitively is not possible with an observational study.⁶ Unmeasured variables, referred to as *lurking variables*, could be responsible for associations observed in the data. By contrast, with an experiment that randomly assigns subjects to treatments, those treatments should balance on any unmeasured variables, at least approximately. For example, in the Harvard Medical School study of the association between heart attack prevalence and taking aspirin or placebo, those taking aspirin would not tend to be

⁴See www.nytimes.com/2019/11/29/business/economics-nobel.html

⁵Such a study is called a *matched case-control study*. It is a *retrospective* type of observational study, which compares groups, such as those with and without some medical condition, by “looking into the past” to measure relevant variables.

⁶Section 6.2.4 discusses this issue in more detail.

younger or of better health than those taking placebo. Because a randomized experiment balances the groups being compared on lurking variables, one can better establish cause and effect with it than with an observational study.

1.4 Descriptive Statistics: Summarizing Data

Descriptive statistics summarize the information that the data contain. Before calculating and analyzing descriptive statistics, you need to be cautious about complications such as missing data or improperly recorded data that can cause software to fail to work or to give invalid output.⁷ Furthermore, the data readily accessible to us may be unstructured and messy. The process of organizing and cleaning the data and bringing them into an appropriate form of a data file for further statistical analysis is called *data wrangling*. This is an important preliminary stage of any statistical analysis and may take substantial time.

This section presents descriptive statistics for quantitative variables in clean data files. Tables and graphs describe the data by showing the number of times various outcomes occurred. The two key features to describe numerically are the *center* of the data and the *variability* of the data around the center.

1.4.1 Example: Carbon Dioxide Emissions in European Nations

Environmental scientists study how the increasing levels of carbon dioxide (CO₂) emissions around the world over time are associated with climate change reflected by the rising temperatures of “global warming.” To illustrate methods of this section, we analyze recent UN data on carbon dioxide emissions per capita, in metric tons, for 31 nations in Europe. The emissions range between 2.0 for Albania and 9.9 for the Netherlands. We do not consider trends over time here, but it has been estimated that CO₂ emissions in the U.S. have doubled since the 1950s and are about 150 times higher than in 1850.

We can use R to read and view the `Carbon` data file from the book’s website:

```
> Carbon <- read.table("http://stat4ds.rwth-aachen.de/data/Carbon.dat", header=TRUE)
      # header=TRUE if variable names are at top of data file
> head(Carbon, 3)      # head(Carbon, n) shows n observations at top of data file
  Nation CO2      # tail(Carbon, n) shows n observations at end of data file
1  Albania  2.0
2  Austria  6.9
3  Belgium  8.3
```

1.4.2 Frequency Distribution and Histogram Graphic

A **frequency distribution** is a listing of the possible values for a variable, together with the number, proportion, or percentage of observations at each value. For a discrete variable with relatively few values, the distribution lists each possible value. For a continuous variable or a discrete variable with many possible values, those values are divided into intervals.

The next output uses R to construct a frequency distribution for the CO₂ values in the `Carbon` data file, using intervals of width 1.0 metric ton. The most common interval was 5.0 up to but not including 6.0, with a percentage occurrence of 25.8%.

⁷Section A.1.4 in the R Appendix shows the use of R software for identifying missing data and conducting descriptive statistical analyses without them.

```

> breaks <- seq(2.0, 10.0, by=1.0) # frequency dist. intervals of width 1 between 2 and 10
> freq <- table(cut(Carbon$CO2, breaks, right=FALSE))
> freq
# with right=FALSE, right-most value not included in interval
[2,3) [3,4) [4,5) [5,6) [6,7) [7,8) [8,9) [9,10)
    1     3     7     8     5     1     3     3
> cbind(freq, freq/nrow(Carbon)) # Frequency distribution of CO2 values, showing
      freq
[2,3) 1 0.03225806
[3,4) 3 0.09677419
[4,5) 7 0.22580645
[5,6) 8 0.25806452
[6,7) 5 0.16129032
[7,8) 1 0.03225806
[8,9) 3 0.09677419
[9,10) 3 0.09677419
> hist(Carbon$CO2, xlab="CO2", ylab="Proportion", freq=FALSE) # histogram
> plot(density(Carbon$CO2)) # smooth-curve approximation of histogram (not shown)

```

A graph of the frequency distribution of a continuous variable or a discrete variable with intervals of values is called a **histogram**.⁸ Each interval of possible values has a bar over it, with height representing its number, proportion or percentage of observations. For the data on European CO2 values, Figure 1.2 shows the histogram constructed using the `hist` function in R. For highly discrete variables, such as categorical variables, a graph that has a separate bar over each distinct value is called a **bar graph**.

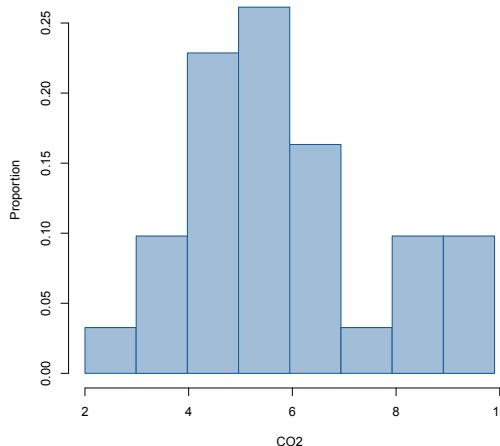


FIGURE 1.2 Histogram for frequency distribution of European CO2 values.

The *shape* of a histogram is informative. For a continuous variable, if we could increase n indefinitely with the number of intervals simultaneously increasing, so their width narrows, the shape would gradually approach a smooth curve.⁹ For a smooth-curve approximation for any n , a bell-shaped appearance indicates that most subjects tend to fall near a central value. See Figure 1.3. The parts of the curve for the lowest values and the highest values are called the *tails* of the distribution. Often, one tail is much longer than the other. A distribution is then said to be **skewed**: *skewed to the right* or *skewed to the left* according to which tail is longer.

⁸Section A.1.2 of the R Appendix presents options for constructing histograms.

⁹This curve can be approximated with the `density` function in R. See the code above (plot not shown) and Exercise 1.18.

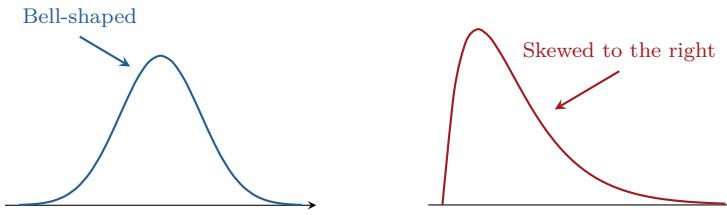


FIGURE 1.3 Smooth curve versions of bell-shaped and skewed frequency distributions: The longer tail indicates the direction of skew.

The **Carbon_West** data file at the book’s website adds four non-European Western nations to the **Carbon** data file for Europe, with CO₂ values of 15.4 for Australia, 15.1 for Canada, 7.7 for New Zealand, and 16.5 for the U.S. As an exercise, load that data file and form the histogram of CO₂ values. The relatively large values for 3 of the 35 nations yields an extended right-tail reflecting skewness to the right.

1.4.3 Describing the Center of the Data: Mean and Median

For a sample, we use subscripts to identify particular observations in a data file. For example, we express the first three observations of a variable denoted by y as y_1, y_2, y_3 .

How can we describe a typical observation for a quantitative variable, for example, describing the *center* of the data? A commonly used measure is the **mean**.

Mean

The **mean** is the sum of the observations divided by the number of them. For a variable y with n observations y_1, y_2, \dots, y_n in a sample from some population, the mean \bar{y} is

$$\bar{y} = \frac{y_1 + y_2 + \dots + y_n}{n} = \frac{\sum_{i=1}^n y_i}{n}.$$

The next section explains how the mean is the *center of gravity* of the data. Because of this, with small n , the mean can be highly influenced by an observation that falls well above or well below the rest of the data. Such an observation is called an **outlier**. Here is an example: The owner of Leonardo’s Pizza Joint reports that the mean annual income of his seven employees is \$56,600. In fact, the annual incomes are \$15,400, \$15,600, \$15,900, \$16,400, \$16,400, \$16,600, and \$299,900. The \$299,900 income, which is the owner’s income, is an outlier. The mean computed for the other six observations alone equals \$16,050, quite different from the overall mean of \$56,600.

The other commonly used measure of center for a quantitative variable is the **median**. It is the middle value,¹⁰ being larger than 50% of the other observations and smaller than 50%. For highly skewed distributions, this is a more representative summary of the center than is the mean. To illustrate, for the ordered income observations for the seven employees of Leonardo’s Pizza, the median is \$16,400. In particular, the median is *resistant* to outliers. When the highest observation of \$299,900 is increased to \$1,000,000, the median is still \$16,400.

A measure sometimes reported with discrete data (quantitative or categorical) is the **mode**, which is the most common outcome. For the GSS data on number of good friends,

¹⁰When n is even, the median is the midpoint between the two middle observations.

the most common response was 2, which is the mode. A frequency distribution that has two distinct peaks is called *bimodal*.

1.4.4 Describing Data Variability: Standard Deviation and Variance

To describe quantitative data more fully, we describe not only the center but also the *variability* about that center. The difference between the largest and smallest observations, called the *range*, is a simple way to do this. The range incorporates only the two most extreme observations, however, so it is sensitive to outliers and does not reflect how far observations other than those two fall from the center.

More useful measures of variability use the distances of *all* the observations from the center. For observation i , the *deviation* of y_i from the mean \bar{y} is $(y_i - \bar{y})$. The deviation is *positive* when y_i falls *above* the mean and *negative* when y_i falls *below* the mean. Calling the mean the *center of gravity* of the data reflects that the sum of the positive deviations equals the negative of the sum of negative deviations. That is, the sum of all the deviations,

$$\sum_{i=1}^n (y_i - \bar{y}) = \sum_{i=1}^n y_i - n\bar{y} = n\bar{y} - n\bar{y} = 0.$$

Because of this, measures of variability use the squares or the absolute values of the deviations.

Standard deviation and variance

For a variable y with n observations y_1, y_2, \dots, y_n in a sample from some population, the *standard deviation* s is

$$s = \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}} = \sqrt{\frac{(y_1 - \bar{y})^2 + (y_2 - \bar{y})^2 + \dots + (y_n - \bar{y})^2}{n-1}}.$$

The standard deviation is the positive square root of the *variance* s^2 ,

$$s^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}.$$

The variance is approximately an average of the squared deviations. The units of measurement are the squares of those for the original data, so the standard deviation is simpler to interpret: s is a sort of *typical distance* of an observation from the mean. Also, *the larger the standard deviation, the greater the spread of the data*. The value of s is nonnegative, with $s = 0$ only when all observations have the same value. The denominator of s^2 uses $(n-1)$, rather than n , because this version of the measure naturally arises in inferential statistical methods.¹¹

The magnitude of s partly reflects the shape of the frequency distribution. If the distribution is approximately bell-shaped, then:¹²

1. About 68% of the observations fall between $\bar{y} - s$ and $\bar{y} + s$.
2. About 95% of the observations fall between $\bar{y} - 2s$ and $\bar{y} + 2s$.

¹¹With data for an entire population, we replace $(n-1)$ by n ; the variance is then precisely the mean squared deviation. In using a sample to *estimate* variability around the population mean, whose value is unknown, Section 4.4.6 shows that the bias due to the numerator of s^2 having the *sample* mean instead of the *population* mean is eliminated when we use $(n-1)$ rather than n in the denominator.

¹²Section 2.5.1 will show where these percentages come from.

3. All or nearly all observations fall between $\bar{y} - 3s$ and $\bar{y} + 3s$.

Figure 1.4 is a graphical portrayal. To illustrate, suppose that grades on an exam are bell-shaped around $\bar{y} = 70$, with $s = 10$. Then, about 68% of the exam scores fall between 60 and 80, about 95% fall between 50 and 90, and all or nearly all fall between 40 and 100.

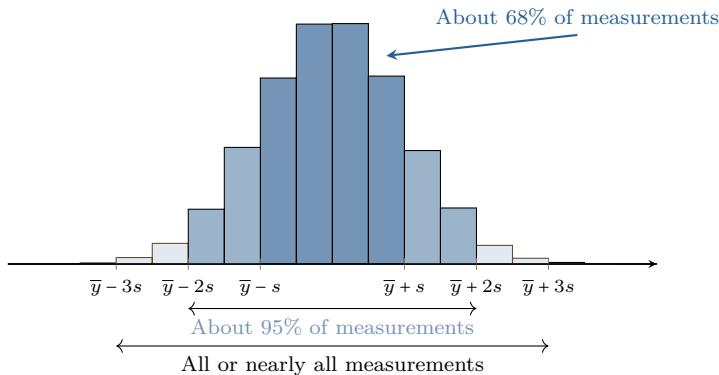


FIGURE 1.4 The standard deviation and mean determine approximate percentages in certain ranges for bell-shaped distributions.

Regardless of the shape of a frequency distribution, it is rare for observations to fall many standard deviations from the mean. The Russian mathematician Pafnuty Chebyshev proved in 1867 that *for any $k \geq 1$, the proportion of observations that fall at least k standard deviations from the mean can be no greater than $1/k^2$* . The result is called **Chebyshev's inequality**.¹³ For example, no more than 4% of the observations can fall at least five standard deviations from the mean. Percentages for most data sets, however, are much closer to the bell-shaped percentages than to the Chebyshev upper bounds.

1.4.5 Describing Position: Percentiles, Quartiles, and Box Plots

Besides center and variability, another way to describe a distribution is with a measure of *position*. The **p th percentile** is the point such that $p\%$ of the observations fall below or at that point and $(100 - p)\%$ fall above it.¹⁴ For example, for $p = 95$, the 95th percentile falls above 95% of the observations and below 5% of them. The 50th percentile is the *median*. **Quantiles** are percentiles expressed in proportion form. For example, the 95th percentile is also called the 0.95 quantile.

Two especially useful percentiles are the the 25th percentile, called the *lower quartile*, and the 75th percentile, called the *upper quartile*. The quartiles together with the median split the distribution into four parts, each containing one-fourth of the observations. The middle half of the data falls between the lower and upper quartiles. The *interquartile range*, denoted by IQR, is the difference between the upper quartile and the lower quartile. Unlike the ordinary range, the IQR is not affected by outliers and takes into account variability by observations other than the most extreme ones.

¹³Sometimes it is referred to as the *Bienaym  -Chebyshev inequality*, since essentially the same result was shown by the French mathematician I. J. Bienaym   in 1853.

¹⁴This definition is imprecise, because no value or an infinite number of real numbers may have *exactly* $p\%$ falling below or at it. Software uses more elaborate definitions that make adjustments to yield precise values. The adjustments are tiny for large n .

The median, the upper and lower quartiles, and the maximum and minimum values provide a ***five-number summary*** of positions. Software can easily find these values as well as other percentiles and summary measures. For instance, with R, here is the five-number summary and the mean and standard deviation for the variable labeled as CO2 in the **Carbon** data file:

```
> summary(Carbon$CO2)           # 1st Qu = lower quartile, 3rd Qu = upper quartile
   Min.   1st Qu.   Median   Mean   3rd Qu.   Max.
   2.000    4.350    5.400   5.819   6.700   9.900
> c(mean(Carbon$CO2), sd(Carbon$CO2), quantile(Carbon$CO2, 0.90))
[1] 5.819355 1.964929 8.900000 # mean, standard deviation, 0.90 quantile
> boxplot(Carbon$CO2, xlab="CO2 values", horizontal=TRUE)
```

The five-number summary is the basis of a graphical display called the ***box plot***. The *box* contains the central 50% of the data, from the lower quartile to the upper quartile. The median is marked by a line drawn within the box. The dashed lines extending from the box, called *whiskers*, contain the outer 50% of the data except for outliers, which are marked separately.¹⁵ For example, Figure 1.5 shows the box plot for the European CO2 values. The shape indicates that the right-tail of the distribution, which corresponds to the relatively large values, is slightly longer than the left tail. The plot reflects the slight skewness to the right of the observations.

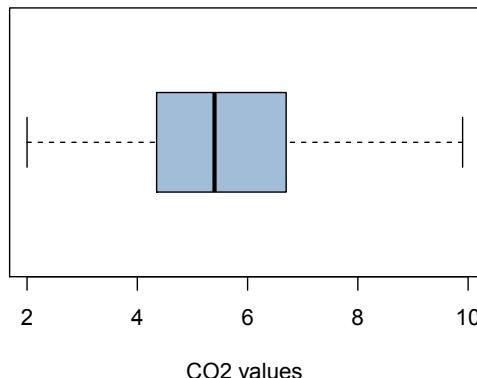


FIGURE 1.5 Box plot of CO2 values for European nations.

Many studies compare different groups on some variable. Side-by-side box plots are useful for making comparisons. To illustrate, Figure 1.6 shows side-by-side box plots of the number of murders per 100,000 population in 2018 for the 50 states and the District of Columbia (D.C.) in the U.S. and for the 10 provinces of Canada. The figure shows that the murder rates in the U.S. tended to be much larger, have much greater variability, and have an extremely large outlier (the murder rate of 24.2 in D.C.). Here is R code to construct this plot and organize summary statistics by group.¹⁶

```
> Crime <- read.table("http://stat4ds.rwth-aachen.de/data/Murder2.dat", header=TRUE)
> boxplot(Crime$murder ~ Crime$nation, xlab="Murder rate", horizontal=TRUE)
```

¹⁵An observation is identified as an *outlier* if it falls more than 1.5(IQR) below the lower quartile or above the upper quartile and as an *extreme outlier* if more than 3(IQR) away. The box plot is one of many methods of ***exploratory data analysis*** proposed in a landmark book by John Tukey (1977).

¹⁶Section A.1.3 shows other ways to present statistics by the level of a second variable.

```
> tapply(Crime$murder, Crime$nation, summary) # applies summary to murder, by nation
$Canada
   Min. 1st Qu. Median   Mean 3rd Qu.   Max.
0.000  1.030  1.735  1.673  1.875  4.070
$US
   Min. 1st Qu. Median   Mean 3rd Qu.   Max.
1.000  2.650  5.000  5.253  6.450 24.200
```

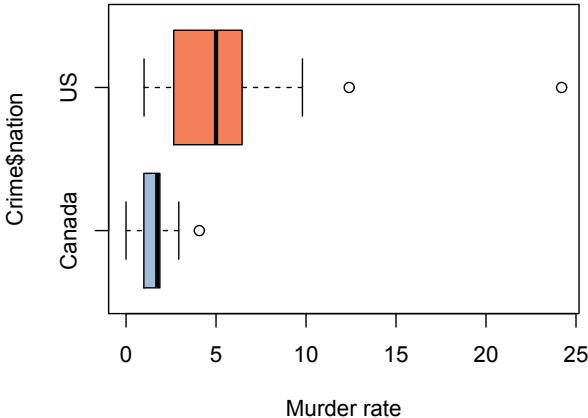


FIGURE 1.6 Side-by-side box plots for U.S. and Canadian murder rates.

1.5 Descriptive Statistics: Summarizing Multivariate Data

For statistical analyses with two or more variables, descriptive methods investigate *associations* between the response variable and the explanatory variables. For a quantitative response variable and categorical explanatory variable, we can compare means or medians on the quantitative variable for the groups that are formed by the categories of the explanatory variable and display results in a side-by-side box plot, such as the R output above and Figure 1.6 show for murder rate and nation. We next present descriptive methods for pairs of quantitative variables and for pairs of categorical variables.

1.5.1 Bivariate Quantitative Data: The Scatterplot, Correlation, and Regression

For a pair of quantitative variables, we can plot values for the explanatory variable on the horizontal (x) axis and for the response variable on the vertical (y) axis. The values of the two variables for any particular observation form a point relative to these axes. A *scatterplot* portrays the n observations as n points.

For example, Figure 1.7 plots the relation in the U.S. for the 50 states and D.C. between x = percent of people in state who own guns and y = suicide rate, measured as the annual number of suicides per 100,000 people in the state. The scatterplot for the 51 observations shows that relatively high values of x tend to occur with relatively high values of y , partly reflecting that guns are used in slightly more than half of all suicides in the U.S.

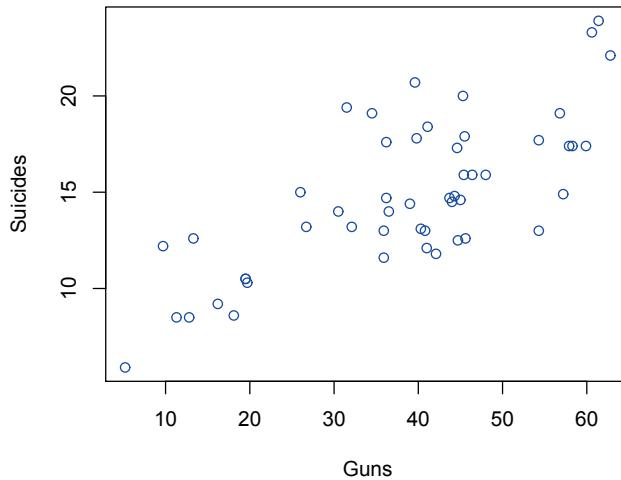


FIGURE 1.7 Scatterplot relating state-level data in the U.S. on percent gun ownership and suicide rate.

Chapter 6 presents two ways to describe such a trend. The **correlation** describes the strength of the association, in terms of how closely the data follow a *straight line trend*. For Figure 1.7, the correlation is 0.74. The positive value means that the suicide rate tends to go *up* as gun ownership goes *up*. The correlation takes values between -1 and $+1$. The larger it is in absolute value, the stronger the association. A **regression analysis** provides a straight-line formula for predicting the response variable from the explanatory variable. For the data in Figure 1.7, this equation is

$$\text{Predicted suicide rate} = 7.390 + 0.1936(\text{gun ownership}).$$

For a state with $x = 5.2$ (the lowest value in this sample, which is for D.C.), the predicted suicide rate is $7.390 + 0.1936(5.2) = 8.40$. For a state with $x = 62.8$ (the highest value, which is for Wyoming), the predicted suicide rate is $7.390 + 0.1936(62.8) = 19.55$.

Chapter 6 derives formulas for the correlation and the regression line. It is simple to implement them with software, as shown next with R.

```
> GS <- read.table("http://stat4ds.rwth-aachen.de/data/Guns_Suicide.dat", header=TRUE)
> Guns <- GS$guns; Suicides <- GS$suicide
> plot(Guns, Suicides)           # scatterplot with arguments x, y
> cor(Guns, Suicides)          # correlation
[1] 0.7386667
> summary(lm(Suicides ~ Guns)) # lm(y ~ x) is "linear model" for
Coefficients:             # response variable y and explanatory variable x
              Estimate
(Intercept) 7.39008      # intercept of line fitted to data points
Guns        0.19356      # slope of line
```

1.5.2 Bivariate Categorical Data: Contingency Tables

For two categorical variables, a **contingency table** is a rectangular table that cross-classifies the variables. The *cells* show the combinations of categories and their counts. For instance,

TABLE 1.1 Contingency table cross-classifying political party identification (ID) and race.

		Political Party ID		
		Democrat	Independent	Republican
Race				
Black		281	66	30
Other		124	77	52
White		633	272	704

Table 1.1 is a contingency table for political party identification (ID = Democrat, Independent, or Republican) and race, using data from the 2018 General Subject Survey.

Treating political party ID as the response variable, we can summarize by finding percentages in each ID category, by race. For instance, 44% of Whites and 8% of Blacks identified as Republicans. The following code uses R with the `PartyID` data file at the book's website to construct the contingency table and find the ID proportions separately for each category of race:

```
> PID <- read.table("http://stat4ds.rwth-aachen.de/data/PartyID.dat", header=TRUE)
> PID
  race      id
1   white    Democrat
...
2238 other    Republican # 2238 subjects in data file
> table(PID$race, PID$id) # forms contingency table (not shown here; see Table 1.1)
> options(digits=2)
> prop.table(table(PID$race, PID$id), margin=1) # For margin=1, proportions
  Democrat Independent Republican # sum to 1.0 within rows
  black     0.75      0.17     0.08
  other     0.49      0.30     0.21
  white     0.39      0.17     0.44
> mosaicplot(table(PartyID$race, PartyID$id)) # graphical portrayal of cell sizes
```

As an exercise, check what the `mosaicplot` function portrays. Contingency tables extend to multi-dimensional tables to handle several variables at once.

1.5.3 Descriptive Statistics for Samples and for Populations

This chapter has introduced commonly-used graphical and numerical descriptive statistics. We introduce others in later chapters and in the book and website appendices. Of descriptive statistics, the mean \bar{y} and the standard deviation s are the most commonly reported measures of center and variability. The correlation and the regression slope are commonly reported measures of association. Since the values of these statistics depend on the sample selected, they vary in value from sample to sample. In this sense, they are also variables.

For example, for a population, Stanford–Binet IQ scores have a bell shape and are scaled to have a mean of 100 and standard deviation of 16. The following R code randomly samples 30 people from a bell-shaped population with this center and spread and finds \bar{y} and s and constructs a histogram. Then it repeats with another sample of size 30, to show that \bar{y} and s and the histogram shape vary from sample to sample.

```
> y1 <- rnorm(30, 100, 16)      # randomly sample normal distribution (Sec. 2.5.1)
> mean(y1); sd(y1); hist(y1)    # histogram (not shown)
[1] 102.7339                   # mean
[1] 11.64643                   # standard deviation
> y2 <- rnorm(30, 100, 16)      # another random sample of size n=30
> mean(y2); sd(y2); hist(y2)
[1] 99.17068      # mean and standard deviation change for each sample of size n=30
[1] 16.52736
```

Do this several times to investigate how \bar{y} and s vary from sample to sample around the population values. You will observe that the histogram may portray the population poorly, sometimes even showing multiple modes. Now do this several times with $n = 1000$ instead of 30. The statistics will vary less from sample to sample, and the histogram will better portray the population. Chapter 3 shows that, for a particular value of n , we can predict how much \bar{y} varies among samples.

Inferential statistical methods use sample descriptive statistics to make predictions about corresponding parameters for the population. In this text, *lower-case Greek letters denote population parameters and Roman letters denote sample statistics*. For example, \bar{y} and s denote a sample mean and standard deviation, and μ and σ denote the population mean and standard deviation. For the IQ sampling just shown, $\mu = 100$ and $\sigma = 16$, whereas $\bar{y} = 102.73$ and $s = 11.65$ for the first sample of 30 observations and $\bar{y} = 99.17$ and $s = 16.53$ for the second sample. We use π for a population proportion.

Before learning about inferential statistical methods, however, you need some basic tools of *probability*, which serves as the language for expressing uncertainty about inferences. Probability is the subject of Chapter 2.

1.6 Chapter Summary

The field of statistical science includes methods for

- designing research studies,
- summarizing the data (*descriptive statistics*),
- making estimations and predictions using the data (*inferential statistics*).

Statistical methods apply to observations in a *sample* taken from a *population*. With randomization, samples are likely to be representative of the population. For a *simple random sample*, every possible sample has the same chance of selection. *Statistics* summarize sample data, while *parameters* summarize entire populations. Inferential statistics use sample data to make predictions about population parameters.

A *data file* has a separate row of data for each subject and a separate column for each characteristic. Statistical methods analyze data on *variables*, which are characteristics that vary among subjects.

- Numerically measured variables, such as family income and number of children in a family, are *quantitative*.
- Variables taking values in a set of categories, such as race and gender, are *categorical*.
- Variables are also classified as *discrete*, such as categorical variables and quantitative variables that take values in a set of separate numbers (e.g., 0, 1, 2, ...), or *continuous*, having a continuous, infinite set of possible values.

Descriptive statistics summarize key characteristics of the data. A *frequency distribution* lists numbers of observations for possible values or intervals of values of a variable. For a quantitative variable, a *histogram* uses bars over possible values or intervals of values to portray a frequency distribution. It shows whether the distribution is approximately bell-shaped or skewed to the right (longer tail pointing to the right) or to the left. A *box*

plot portrays the quartiles (25th and 75th percentiles), the extreme values, and any outliers. **Table 1.2** summarizes the most important numerical measures that describe the *center* of n observations and their *variability* (spread).

TABLE 1.2 Summary of descriptive statistical measures of center and variability.

Measure	Definition	Interpretation
Center		
Mean	$\bar{y} = (\sum_i y_i)/n$	Center of gravity
Median	Middle observation of ordered sample	50th percentile, splits sample into two equal parts
Variability		
Standard deviation	$s = \sqrt{\sum(y_i - \bar{y})^2/(n - 1)}$	If bell-shaped, 68%, 95%, nearly all within $s, 2s, 3s$ of \bar{y}

Bivariate statistics summarize the *association* between two variables and how the outcome on a *response variable* depends on the value of an *explanatory variable*.

- For quantitative variables, a *scatterplot* graphs the observations as points with axes for the variables. The *correlation* describes the strength of straight-line association and the *regression line* can predict the response variable using the explanatory variable.
- For categorical variables, a *contingency table* shows the number of observations at the combinations of possible category outcomes for the two variables.

Exercises

Data Analysis and Applications

- 1.1 In the 2018 election for Senate in California, a CNN exit poll of 1882 voters stated that 52.5% voted for the Democratic candidate, Diane Feinstein. Of all 11.1 million voters, 54.2% voted for Feinstein.
- What was the (i) subject, (ii) sample, (iii) population?
 - Identify a relevant statistic and corresponding parameter.
- 1.2 The **Students** data file at <http://stat4ds.rwth-aachen.de/data> shows responses of a class of 60 social science graduate students at the University of Florida to a questionnaire that asked about *gender* (1 = female, 0 = male), *age*, *hsgpa* = high school GPA (on a four-point scale), *cogpa* = college GPA, *dhome* = distance (in miles) of the campus from your home town, *dres* = distance (in miles) of the classroom from your current residence, *tv* = average number of hours per week that you watch TV, *sport* = average number of hours per week that you participate in sports or have other physical exercise, *news* = number of times a week you read a newspaper, *aids* = number of people you know who have died from AIDS or who are HIV+, *veg* = whether you are a vegetarian (1 = yes, 0 = no), *affil* = political affiliation (1 = Democrat, 2 = Republican, 3 = independent), *ideol* = political ideology (1 = very liberal, 2 = liberal, 3 = slightly liberal, 4 = moderate, 5 = slightly conservative, 6 = conservative, 7 = very conservative), *relig* = how often you attend religious services (0 = never, 1 = occasionally, 2 = most weeks, 3 = every

week), *abor* = opinion about whether abortion should be legal in the first three months of pregnancy (1 = yes, 0 = no), *affirm* = support affirmative action (1 = yes, 0 = no), and *life* = belief in life after death (1 = yes, 2 = no, 3 = undecided). You will use this data file for some exercises in this book.

- (a) Practice accessing a data file for statistical analysis with your software by going to the book's website and copying and then displaying this data file.
 - (b) Using responses on *abor*, state a question that could be addressed with (i) descriptive statistics, (ii) inferential statistics.
- 1.3 Identify each of the following variables as categorical or quantitative: (a) Number of smartphones that you own; (b) County of residence; (c) Choice of diet (vegetarian, nonvegetarian); (d) Distance, in kilometers, commute to work
- 1.4 Give an example of a variable that is (a) categorical; (b) quantitative; (c) discrete; (d) continuous.
- 1.5 In analyzing data about patients who developed Covid-19 from coronavirus, many research studies used the scale (1. Death; 2. Hospitalized with invasive ventilation; 3. Hospitalized with non-invasive ventilation; 4. Hospitalized with supplemental oxygen; 5. Hospitalized, not requiring supplemental oxygen but requiring ongoing medical care; 6. Hospitalized, not requiring ongoing medical care (quarantine or awaiting rehab); 7. Not hospitalized, limitation on activities; 8. Not hospitalized, no limitations on activities). Is this categorical scale *nominal* or *ordinal*? Why?
- 1.6 Give an example of a variable that is (a) technically discrete but essentially continuous for purposes of data analysis; (b) potentially continuous but highly discrete in the way it is measured in practice.
- 1.7 The student directory for a large university has 400 pages with 130 names per page, a total of 52,000 names. Using software, show how to select a simple random sample of 10 names.
- 1.8 Explain whether an experiment or an observational study would be more appropriate to investigate the following:
- (a) Whether cities with higher unemployment rates tend to have higher crime rates.
 - (b) Whether a Honda Accord hybrid or a Toyota Prius gets better gas mileage.
 - (c) Whether higher college grade point averages tend to occur for students who had higher scores on college entrance exams.
 - (d) Whether design A or design B for an Internet page makes a person more likely to buy the product advertised.
- 1.9 For the GSS data on number of good friends (search at <https://sda.berkeley.edu/sdaweb/analysis/?dataset=gss18> for variable NUMFREN), the responses (1, 2, 3, 4, 5, 6, 7, 8, 9, ..., 96+) had percentages of (6.1, 16.2, 15.7, 14.2, 11.3, 8.8, 1.2, 2.4, 0.7, ..., 0.8). Report the median and the mode. Would you expect the mean to be smaller, the same, or larger than the median? Why?
- 1.10 Analyze the **Carbon_West** data file at the book's website by (a) constructing a frequency distribution and a histogram, (b) finding the mean, median, and standard deviation. Interpret each.

- 1.11 According to Statistics Canada, for the Canadian population having income in 2019, annual income had a median of \$35,000 and mean of \$46,700. What would you predict about the shape of the distribution? Why?
- 1.12 Give an example of a variable that is nonnegative but has the majority of its sample values at 0, so the median is not especially informative despite the skew.
- 1.13 A report indicates that public school teacher's annual salaries in New York city have an approximate mean of \$69,000 and standard deviation of \$6,000. If the distribution has approximately a bell shape, report intervals that contain about (a) 68%, (b) 95%, (c) all or nearly all salaries. Would a salary of \$100,000 be unusual? Why?
- 1.14 In 2017, according to the Kaiser Family Foundation (www.kff.org), the five-number summary for the U.S. statewide percentage of people without health insurance had minimum = 3% (Massachusetts), lower quartile = 6%, median = 8%, upper quartile = 9.5%, and maximum = 17% (Texas). Interpret the quartiles and the interquartile range, and sketch a box plot.
- 1.15 According to www.salary.com, the mean salary (in dollars) of secondary school teachers in the United States in 2019 varied among states with a five-number summary of maximum = 67,600 (California), upper quartile = 64,700, median = 55,500, lower quartile = 53,100, and minimum = 51,600 (South Dakota). Sketch a box plot, and indicate whether the distribution seems to be symmetric, skewed to the right, or skewed to the left.
- 1.16 Access the most recent General Social Survey at <https://sda.berkeley.edu/archive.htm>. Entering TVHOURS for the row variable and year(2018) in the selection filter, you obtain data on hours per day of TV watching in the U.S. in 2018.
- Construct the frequency distribution for the values 0, 1, 2, 3, 4, 5, 6, 7 or more. How would you describe its shape?
 - Find the median and the mode.
 - Check *Summary statistics* in the output options, and report the mean and standard deviation. From these, explain why you would not expect the distribution to be bell-shaped.
- 1.17 From the *Murder* data file at the book's website, use the variable *murder*, which is the murder rate (per 100,000 population) for each state in the U.S. in 2017 according to the FBI Uniform Crime Reports. At first, do not use the observation for D.C. (DC). Using software:
- Find the mean and standard deviation and interpret their values.
 - Find the five-number summary, and construct the corresponding box plot. Interpret.
 - Now include the observation for D.C. What is affected more by this outlier: The mean or the median? The range or the inter-quartile range?
- 1.18 The *Income* data file at the book's website reports annual income values in the U.S., in thousands of dollars.
- Using software, construct a histogram. Describe its shape.
 - Find descriptive statistics to summarize the data. Interpret them.

- (c) The *kernel density estimation* method finds a smooth-curve approximation for a histogram. At each value, it takes into account how many observations are nearby and their distance, with more weight given those closer. Increasing the *bandwidth* increases the influence of observations further away. Plot a smooth-curve approximation for the histogram of income values, using the `density` function in R. Summarize the impact of increasing and of decreasing the bandwidth (option `bw` in the `density` function) substantially from the default value.¹⁷
- (d) Construct and interpret side-by-side box plots of income by race (B = Black, H = Hispanic, W = White). Compare the incomes using numerical descriptive statistics.
- 1.19 The **Houses** data file at the book's website lists the selling price (thousands of dollars), size (square feet), tax bill (dollars), number of bathrooms, number of bedrooms, and whether the house is new (1 = yes, 0 = no) for 100 home sales in Gainesville, Florida. Let's analyze the selling prices.
- Construct a frequency distribution and a histogram. Describe the shape.
 - Find the percentage of observations that fall within one standard deviation of the mean. Why is this not close to 68%?
 - Construct a box plot, and interpret.
 - Use descriptive statistics to compare selling prices according to whether the house is new.
- 1.20 Refer to the previous exercise. Let y = selling price and x = size of home.
- Construct a scatterplot. Interpret. Identify any observation that seems to fall apart from the others.
 - Find the correlation. Interpret.
 - Find the regression line. Interpret the slope, and find the predicted selling price for a home of (i) 1000 square feet, (ii) 4000 square feet.
- 1.21 For the **Students** data file introduced in Exercise 1.2, summarize the relationship between *hsgpa* and *cogpa* using correlation and regression. Find the predicted college GPA of a student who had a high school GPA of 4.0.
- 1.22 Using the **Happy** data file, construct the contingency table relating marital status and happiness. Which variable is the natural response variable? Report the proportions in its categories, separately at each category of the explanatory variable, and interpret.
- 1.23 For the **Students** data file introduced in Exercise 1.2, construct and summarize a contingency table relating religiosity and opinion about legalized abortion.
- 1.24 The **UN** data file at the book's website has United Nations data for 42 nations on per capita gross domestic product (GDP, in thousands of dollars), a human development index (HDI, which has components referring to life expectancy at birth, educational attainment, and income per capita), a gender inequality index (GII, a composite measure reflecting inequality in achievement between women and men in reproductive health, empowerment, and the labor market), fertility rate (number of births per woman), carbon dioxide emissions per capita (CO2, in metric tons), a homicide rate (number of homicides per 100,000 people), prison population (per 100,000 people), and percent using the Internet.

¹⁷The default is the so-called *Silverman's rule of thumb*, option `bw = "nrd0"`.

- (a) Conduct a descriptive statistical analysis of the prison rates. Summarize your conclusions, highlighting any unusual observations.
- (b) Using a command like `cor(cbind(GDP, HDI, GII, ..., Internet))` in R, construct a *correlation matrix* showing the correlation for each pair of variables. Which pair has the strongest association?
- (c) Conduct correlation and regression analyses to study the association between CO2 use and GDP. Show how the predicted CO2 varies as GDP goes from its minimum to its maximum value in this sample.
- 1.25 The **ScotsRaces** data file at the book's website shows the record times for men and for women for several hill races in Scotland. Use graphical and numerical descriptive statistics to compare the men's winning times with the women's winning times. Summarize your analyses in a short report, with software output as an appendix.
- 1.26 Refer to the previous exercise. Explanatory variables listed in the data file are the distance of the race and the climb in elevation. Use graphical and numerical descriptive statistics to summarize the men's winning times and their relationship with the race distance and climb. Summarize your analyses in a one-page report.
- 1.27 A study of sheep¹⁸ analyzed whether the sheep survived for a year from the original observation time (1 = yes, 0 = no) as a function of their weight (*kg*) at the original observation. Using **Sheep** data file at the text website, use graphical and numerical methods of this chapter to compare weights of the sheep that survived to weights of the sheep that did not survive. Summarize the results of your analysis in a few sentences.

Methods and Concepts

- 1.28 The beginning of [Section 1.2](#) mentioned a potential survey to observe characteristics such as opinion about the legalization of same-sex marriage, political party affiliation, frequency of attending religious services, number of years of education, annual income, marital status, race, and gender. Describe two ways you could select one of these variables as a response variable and the others as explanatory variables, and explain the reasoning for these choices.
- 1.29 A *systematic random sample* of n subjects from a population of size N selects a subject at random from the first $k = N/n$ in the population list and then selects every k th subject listed after that one. Explain why this is not a simple random sample.

For the following two multiple-choice questions, select the best response.

- 1.30 A simple random sample of size n is one in which:
- (a) Every n th member is selected from the population.
- (b) Each possible sample of size n has the same chance of being selected.
- (c) There must be exactly the same proportion of women in the sample as is in the population.
- (d) You keep sampling until you have a fixed number of people having various characteristics (e.g., males, females).

¹⁸Summarized in article by T. Coulson, *Oikos*, **121**: 1337–1350 (2012); thanks to Prof. M. K. Oli for the data from this study.

- (e) A particular minority group member of the population is less likely to be chosen than a particular majority group member.
- (f) All of the above.
- 1.31 If we use random numbers to take a simple random sample of 50 students from the 6500 undergraduate students at the University of Rochester:
- We would never get the random number 1111, because it is not a random sequence.
 - The draw 1234 is no more or less likely than the draw 1111.
 - Since the sample is random, it is *impossible* that it will be non-representative, such as having only females in the sample.
 - Since the sample is random, it is impossible to get the sequence of random numbers 0001, 0002, 0003, ..., 0049, 0050.
- 1.32 With an Internet search, find a study that used an **(a)** experiment; **(b)** observational study. In each case, describe how the sample was obtained and summarize results.
- 1.33 An article¹⁹ in the *New England Journal of Medicine* (October 12, 2012) observed a correlation of 0.79 for 23 countries between per capita annual chocolate consumption and the number of Nobel laureates per 10 million population. Was this study an experiment or an observational study? Can we conclude that increasing chocolate consumption increases the chance of a Nobel prize? Why or why not?
- 1.34 A research study funded by Wobegon Springs Mineral Water, Inc., discovers that children with dental problems are less common in families that regularly buy bottled water than in families that do not. Explain why this association need not reflect a causal link between drinking bottled water and having fewer dental problems. Identify lurking variables that could be responsible for the result, and explain how.
- 1.35 Suppose that grade-point averages at your university are bell-shaped with mean 3.0 and standard deviation 0.3. Randomly sample n students several times for $n = 20$ (using a function such as `rnorm` in R), and then repeat several times for $n = 1000$, each time constructing a histogram. What does this suggest about the difficulty of determining the shape of a population distribution when n is small?
- 1.36 Construct a set of data for which the mean and median are identical.
- 1.37 Suppose you estimate the mean number of friends that members of Facebook have by randomly sampling 100 members of Facebook and (i) averaging the numbers of friends that they have, (ii) averaging how many friends the friends of those members have. Which estimate do you think would be larger? Why? (*Hint:* See article by J. A. Paulos, *Scientific American*, Feb. 1, 2011.)
- 1.38 To measure center, why is the **(a)** median sometimes preferred over the mean? **(b)** mean sometimes preferred over the median? (*Hint:* A wide variety of highly discrete frequency distributions can have the same median.) To illustrate, in your answers use the variables annual income for (a) and number of times you played a sport in the past week for (b).
- 1.39 To measure variability, why is the **(a)** standard deviation s usually preferred over the range? **(b)** interquartile range often preferred over the range?

¹⁹See www.nejm.org/doi/full/10.1056/NEJMoa1211064

- 1.40 The largest value in a sample is moved upwards so that it is an extreme outlier. Explain how, if at all, this affects the mean, median, range, and interquartile range.
- 1.41 To investigate how \bar{y} can vary from sample to sample of size n , for the simulation from a bell-shaped population shown at the end of [Section 1.5.3](#), take (a) 10,000 random samples of size $n = 30$ each; (b) 10,000 random samples of size $n = 1000$ each. In each case, form a histogram of the 10,000 \bar{y} values and find their standard deviation. Compare results and explain what this simulation reveals about the impact of sample size on how study results can vary. ([Chapter 3](#) shows that in sampling from a population with standard deviation 16, the theoretical standard deviation of \bar{y} values is $16/\sqrt{n}$.)
- 1.42 The Internet site www.artofstat.com/web-apps has useful apps²⁰ for illustrating data analyses and properties of statistical methods.
- (a) Using the *Explore Quantitative Data* app, construct a sample of 20 observations on y = number of hours of physical exercise in the past week having $\bar{y} < s$. What aspect of the distribution causes this to happen?
 - (b) Using the *Explore Linear Regression* app with the *Draw Own* option, create 20 data points that are plausible for x = number of hours of exercise last week and y = number of hours of exercise this week. Describe your data by the correlation and by the linear regression line, and interpret them.
- 1.43 For a sample with mean \bar{y} , show that adding a constant c to each observation changes the mean to $\bar{y} + c$, and the standard deviation s is unchanged. Show that multiplying each observation by c changes the mean to $c\bar{y}$ and the standard deviation to $|c|s$.
- 1.44 Suppose the sample data distribution of $\{y_i\} = \{y_1, \dots, y_n\}$ is very highly skewed to the right, and we take logs and analyze $\{x_i = \log(y_i)\}$.
- (a) Is $\bar{x} = \log(\bar{y})$? Why or why not?
 - (b) Is $\text{median}(\{x_i\}) = \log[\text{median}(\{y_i\})]$? Why or why not?
 - (c) To summarize $\{y_i\}$, we find \bar{x} and then use $\exp(\bar{x})$. Show that $\exp(\bar{x}) = (\prod_i y_i)^{1/n}$, called the ***geometric mean*** of $\{y_i\}$.
- 1.45 Find the Chebyshev inequality upper bound for the proportion of observations falling at least (a) 1, (b) 2, (c) 3 standard deviations from the mean. Compare this to the approximate proportions for a bell-shaped distribution. Why are the differences so large?
- 1.46 The R output in [Section 1.5.2](#) shows a `mosaicplot` function. Implement it for the `PartyID` data. Do an Internet search and look at the R Appendix to this book to learn about *mosaic plots*, and describe what the plot shows you. (The `mosaic` function in the `vcd` package constructs more sophisticated mosaic plots.)
- 1.47 The *least squares* property of the mean states that the data fall closer to \bar{y} than to any other number c , in the sense that

$$\sum_i (y_i - \bar{y})^2 < \sum_i (y_i - c)^2.$$

Prove this property by treating $f(c) = \sum_i (y_i - c)^2$ as a function of c and deriving the value of c that minimizes it.

²⁰The apps at that website were developed by Dr. Bernhard Klingenberg.

- 1.48 For a sample $\{y_i\}$ of size n , $\sum_i |y_i - c|$ is minimized at $c = \text{median}$. Explain why this property holds. (Hint: Starting at $c = \text{median}$, what happens to $\sum_i |y_i - c|$ as you move away from it in either direction?)
- 1.49 The Delphi group at Carnegie-Mellon University has tracked statistics about coronavirus. [Figure 1.8](#) is a scatterplot of U.S. statewide data compiled between x = percentage wearing masks and y = percentage knowing someone with Covid-19 symptoms. Which do you think best describes the scatterplot?

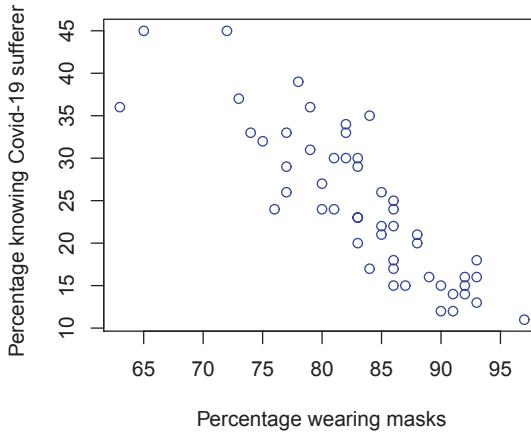


FIGURE 1.8 Scatterplot for data file `CovidMasks` at book's website, showing U.S. statewide data in October 2020 of percentage of people wearing masks in public all or most of the time and percentage of people knowing someone who has had Covid-19 symptoms

- (a) Correlation = -0.85 and predicted $y = 112.9 - 1.06 x$
- (b) Correlation = -0.85 and predicted $y = 112.9 + 1.06 x$
- (c) Correlation = 0.85 and predicted $y = 112.9 - 1.06 x$
- (d) Correlation = 0.85 and predicted $y = 112.9 + 1.06 x$

2

Probability Distributions

This chapter introduces *probability*, the language used to describe uncertainty, such as in inferential statistical analyses. We first define probability and give basic rules for calculating probabilities. *Random variables* specify possible values of variables in experiments or observational studies that incorporate randomization. A *probability distribution* summarizes random variation in the values of a random variable by specifying probabilities for all the possible outcomes. Probability distributions are themselves summarized by measures of center and variability, such as the mean and standard deviation. These measures are *expected values*, describing what we expect, on the average, with observations taken at random from a probability distribution.

Like ordinary variables, random variables can be *discrete* or *continuous*. The most important probability distributions for discrete random variables are the *binomial* for binary outcomes and the *Poisson* for count outcomes. The most important probability distributions for continuous random variables are the *normal*, which has a bell shape over the real line, and the *gamma*, which is skewed to the right over positive real number values.

Joint probability distributions describe how two or more random variables co-vary. The *correlation* describes their strength of association, with *independence* being the absence of any association.

2.1 Introduction to Probability

In everyday life, the term *probability* is used in an informal manner to mean the *chance* of an outcome. But this definition is rather vague. For studies that can employ *randomization* for gathering data, a more precise definition relates to how often that outcome occurs.

2.1.1 Probabilities and Long-Run Relative Frequencies

For random phenomena, such as observations in a randomized experiment, the probability of an outcome refers to the proportion of times the outcome would occur in a very long sequence of like observations.

Probability

For an observation of a random phenomenon, the **probability** of a particular outcome is the proportion of times that outcome would occur in an indefinitely long sequence of like observations, under the same conditions.

Probabilities are sometimes expressed as percentages. For this definition, a weather forecast that the probability of rain tomorrow is 0.20, or 20%, means that in a long run of repeated observations of days with atmospheric conditions like those expected tomorrow, rain would occur on 20% of the days. We could simulate the outcome for tomorrow by

generating a random digit from $(0, 1, 2, \dots, 9)$, where 2 of the 10 possible outcomes (say, 0 and 1) denote *rain* and the other 8 possible outcomes denote *no rain*; for example, using R:

```
> sample(0:9, 1)      # randomly generate 1 integer between 0 and 9
[1] 8
```

The outcome of the simulation was an 8, representing no rain. Let's simulate for a week under the same conditions:

```
> sample(0:9, 7, replace=TRUE)    # with replace=TRUE, numbers replaced,
[1] 5 8 1 9 8 9 5                # each selection has same choice set
```

In this simulation, it rains on day 3 but not on any of the other six days.

The outcome on a particular day corresponds to flipping an unbalanced coin that has probability 0.20 of a head, where we identify the *head* outcome with *rain*. More generally, we can simulate how often rain occurs in n days by flipping a coin n times, when each flip has probability 0.20 of a head and probability 0.80 of a tail. We can do this using the R function `rbinom`, which provides simulations with binary outcomes for each observation.¹ Here is the result for 1 simulation of $n = 7$ days:

```
> rbinom(1, 7, 0.20) # 1 simulation of 7 coin flips, probability 0.20 of head
[1] 3                  # outcome is 3 heads in 7 flips (i.e., rain on 3 days)
> rbinom(7, 1, 0.20) # Or, 7 simulations of 1 coin flip, probability 0.20 of head
[1] 0 1 0 0 0 0 1      # 1 = head, 0 = tail; simulated outcome is rain on days 2 and 7
```

For the first simulation, it rained on 3 of the 7 days. The proportion $3/7 = 0.43$ is quite far from the stated probability of 0.20. But a probability is a *long-run* relative frequency, in theory letting the number of days $n \rightarrow \infty$. Let's see what happens with much larger n for a simulation:

```
> rbinom(1, 100, 0.20)
[1] 18                  # proportion 0.180; 18 heads in 100 coin flips
> rbinom(1, 1000, 0.20)
[1] 204                 # proportion 0.204
> rbinom(1, 10000, 0.20)
[1] 2010                # proportion 0.2010
> rbinom(1, 100000, 0.20)
[1] 20032               # proportion 0.2003
> rbinom(1, 1000000, 0.20)
[1] 199859              # proportion 0.1999; 199859 heads in 1000000 coin flips
```

According to this simulation, the proportions of days of rain out of $n = (100, 1000, 10000, 100000, 1000000)$ days are $(0.180, 0.204, 0.201, 0.2003, 0.1999)$. The proportion converges to 0.2000 as n increases indefinitely. Figure 2.1 illustrates how the proportion stabilizes as n increases.

The long-run sequence of like observations for a definition of probability is not always appropriate. For example, it is not meaningful for the probability that intelligent life exists elsewhere in the universe or for the probability that a new business is successful, because no long-run sequence of observations is available. We must then rely on *subjective* beliefs rather than *objective* data. In fact, an alternative definition of probability is subjective:

Subjective definition of probability: The probability of an outcome is the degree of belief that the outcome will occur, based on all the available information.

¹Besides assuming the same probability for each coin flip, this function assumes that observations are “independent events” in a sense to be introduced in Section 2.1.6; for instance, whether it rains one day is not influenced by whether it rained the previous day.